

# Appendix for DeltaGrad: Rapid retraining of machine learning models

Yinjun Wu, Edgar Dobriban, Susan B. Davidson

## Contents

|          |  |           |
|----------|--|-----------|
| <b>A</b> | <b>Mathematical details</b>  | <b>2</b>  |
| A.1      | Additional notes on setup, preliminaries . . . . .                                       | 2         |
| A.1.1    | Classical results on GD convergence, SGD convergence . . . . .                           | 2         |
| A.1.2    | Notations for DeltaGrad with SGD . . . . .   | 3         |
| A.1.3    | Classical results for random variables . . . . .   | 3         |
| A.2      | Results for deterministic gradient descent . . . . .                                     | 4         |
| A.2.1    | Quasi-Newton . . . . .   | 4         |
| A.2.2    | Proof that Quasi-Hessians are well-conditioned . . . . .                                 | 5         |
| A.2.3    | Proof preliminaries . . . . .  | 6         |
| A.2.4    | Main recursions . . . . .  | 10        |
| A.2.5    | Proof of Theorem 2 . . . . .   | 11        |
| A.2.6    | Proof of Theorem 3 . . . . .   | 12        |
| A.2.7    | Proof of Theorem 4 . . . . .   | 13        |
| A.2.8    | Proof of Theorem 5 . . . . .   | 16        |
| A.3      | Results for stochastic gradient descent . . . . .  | 20        |
| A.3.1    | Quasi-Newton . . . . .   | 20        |
| A.3.2    | Proof preliminaries . . . . .  | 20        |
| A.3.3    | Main recursions . . . . .  | 26        |
| A.3.4    | Proof of Theorem 8 . . . . .   | 27        |
| A.3.5    | Proof of Theorem 9 . . . . .   | 28        |
| A.3.6    | Proof of Theorem 10 . . . . .  | 29        |
| A.3.7    | Proof of Theorem 11 . . . . .  | 32        |
| <b>B</b> | <b>Details on applications</b>   | <b>33</b> |
| B.1      | Privacy related data deletion . . . . .  | 33        |
| <b>C</b> | <b>Supplementary algorithm details</b>   | <b>35</b> |
| C.1      | Extension of DeltaGrad for stochastic gradient descent . . . . .                         | 35        |
| C.2      | Extension of DeltaGrad for online deletion/addition . . . . .                            | 35        |
| C.2.1    | Convergence rate analysis for online gradient descent version of DeltaGrad . . . . .     | 35        |
| C.3      | Extension of DeltaGrad for non-strongly convex, non-smooth objective functions . . . . . | 47        |
| <b>D</b> | <b>Supplementary experiments</b>   | <b>47</b> |
| D.1      | Experiments with large deletion rate . . . . .   | 47        |
| D.2      | Influence of hyper-parameters on performance . . . . .                                   | 49        |
| D.3      | Comparison against the state-of-the-art work . . . . .                                   | 51        |
| D.4      | Experiments on large ML models . . . . .   | 51        |
| D.5      | Applications of DeltaGrad to robust learning . . . . .                                   | 52        |

## A Mathematical details

The main result for DeltaGrad with GD is Theorem 5, proved in Section A.2.8.

### A.1 Additional notes on setup, preliminaries

#### A.1.1 Classical results on GD convergence, SGD convergence

**Lemma 1** (GD convergence, folklore, e.g., Boyd and Vandenberghe (2004)). *Gradient descent over a strongly convex objective function with fixed step size  $\eta_t = \eta \leq \frac{2}{L+\mu}$  has exponential convergence rate, i.e.:*

$$F(\mathbf{w}_t) - F(\mathbf{w}^*) \leq c^t \frac{L}{2} \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \quad (1)$$

where  $c := (L - \mu)/(L + \mu) < 1$ .

Recall also that the eigenvalues of the "contraction operator"  $\mathbf{I} - \eta_t \mathbf{H}(\mathbf{w})$  are bounded as follows.

**Lemma 2** (Classical bound on eigenvalues of the "contraction operator"). *Under the convergence conditions of gradient descent with fixed step size, i.e.  $\eta_t = \eta \leq \frac{2}{\mu+L}$ , the following inequality holds for any parameter  $\mathbf{w}$ :*

$$\|\mathbf{I} - \eta \mathbf{H}(\mathbf{w})\| \leq 1. \quad (2)$$

This lemma follows directly, because the eigenvalues of  $\mathbf{I} - \eta \mathbf{H}$  are bounded between  $-1 \leq 1 - \eta L \leq 1 - \eta \mu \leq 1$ .

**Lemma 3** (SGD convergence, see e.g., Bottou et al. (2018)). *Suppose that the stochastic gradient estimates are correlated with the true gradient, and bounded in the following way. There exist two scalars  $J_1 \geq J_2 > 0$  such that for arbitrary  $\mathcal{B}_t$ , the following two inequalities hold:*

$$\nabla F(\mathbf{w}_t)^T \mathbb{E} \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t) \geq J_2 \|\nabla F(\mathbf{w}_t)\|^2, \quad (3)$$

$$\|\mathbb{E} \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t)\| \leq J_1 \|\nabla F(\mathbf{w}_t)\|.$$

Also, assume that for two scalars  $J_3, J_4 \geq 0$ , we have:

$$\text{Var} \left( \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t) \right) \leq J_3 + J_4 \|\nabla F(\mathbf{w}_t)\|^2. \quad (4)$$

By combining equations (3)-(4), the following inequality holds:

$$\mathbb{E} \left\| \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t) \right\|^2 \leq J_3 + J_5 \|\nabla F(\mathbf{w}_t)\|^2$$

where  $J_5 = J_4 + J_1^2 \geq J_2^2 \geq 0$ .

Then stochastic gradient descent with fixed step size  $\eta_t = \eta \leq \frac{J_2}{LJ_5}$  has the convergence rate:

$$\mathbb{E} [F(\mathbf{w}_t) - F(\mathbf{w}^*)] \leq \frac{\eta L J_3}{2\mu J_2} + (1 - \eta \mu J_2)^{t-1} \left( F(\mathbf{w}_1) - F(\mathbf{w}^*) - \frac{\eta L J_3}{2\mu J_2} \right) \rightarrow \frac{\eta L J_3}{2\mu J_2}.$$

If the gradient estimates are unbiased, then  $\mathbb{E} \frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t) = \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}_t) = \nabla F(\mathbf{w}_t)$  and thus  $J_1 = J_2 = 1$ . Moreover,  $J_3 \sim 1/m$ , where  $m$  is the minibatch size, because  $J_2$  is the variance of the stochastic gradient.

So the convergence condition for fixed step size becomes  $\eta_t = \eta \leq \frac{1}{LJ_5}$ , in which  $J_5 = J_4 + J_1^2 = J_4 + 1 \geq 1$ . So  $\eta_t = \eta \leq \frac{1}{LJ_5} \leq \frac{1}{L}$  suffices to ensure convergence.

### A.1.2 Notations for DeltaGrad with SGD

The SGD parameters trained over the full dataset, explicitly trained over the remaining dataset and incrementally trained over the remaining dataset are denoted by  $\mathbf{w}^S$ ,  $\mathbf{w}^{U,S}$  and  $\mathbf{w}^{I,S}$  respectively. Then given the mini-batch size  $B$ , mini-batch  $\mathcal{B}_t$ , the number of removed samples from each mini-batch  $\Delta B_t$  and the set of removed samples  $R$ , the update rules for the three parameters are:

$$\mathbf{w}^S_{t+1} = \mathbf{w}^S_t - \eta \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_t) = \mathbf{w}^S_t - \eta G_{B,S}(\mathbf{w}^S_t), \quad (5)$$

$$\begin{aligned} \mathbf{w}^{U,S}_{t+1} &= \mathbf{w}^{U,S}_t - \eta \frac{1}{B - \Delta B_t} \sum_{i \in \mathcal{B}_t, i \notin R} \nabla F_i(\mathbf{w}^{U,S}_t) \\ &= \mathbf{w}^{U,S}_t - \eta G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t), \end{aligned} \quad (6)$$

$$\mathbf{w}^{I,S}_{t+1} = \begin{cases} \mathbf{w}^{I,S}_t - \frac{\eta}{B-\Delta B_t} \sum_{i \in \mathcal{B}_t, i \notin R} \nabla F(\mathbf{w}^{I,S}_t) & (t - j_0) \bmod T_0 = 0 \\ \mathbf{w}^{I,S}_t - \frac{\eta}{B-\Delta B_t} \{B[\mathbf{B}_{j_m}(\mathbf{w}^{I,S}_t - \mathbf{w}^S_t) + \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_t)] - \sum_{i \in R, i \in \mathcal{B}_t} \nabla F(\mathbf{w}^{I,S}_t)\} & \text{or } t \leq j_0 \\ & \text{otherwise} \end{cases} \quad (7)$$

in which  $G_{B,S}(\mathbf{w}^S_t)$  and  $G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t)$  represent the average gradients over the minibatch  $\mathcal{B}_t$  before and after removing samples.

We assume that the minibatch randomness of  $\mathbf{w}^{U,S}$  and  $\mathbf{w}^{I,S}$  is the same as  $\mathbf{w}^S$ . By following Lemma 3, we assume that the gradient estimates of SGD are unbiased, i.e.  $\mathbb{E}\left(\frac{1}{B_t} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w})\right) = \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}) = \nabla F(\mathbf{w})$  for any  $\mathbf{w}$ , which indicates that:

$$\begin{aligned} \mathbb{E}\left(\frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_t)\right) &= \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}^S_t) = \nabla F(\mathbf{w}^S_t), \\ \mathbb{E}\left(\frac{1}{B - \Delta B_t} \sum_{i \in \mathcal{B}_t, i \notin R} \nabla F_i(\mathbf{w}^{U,S}_t)\right) &= \frac{1}{n - \Delta n} \sum_{i \notin R} \nabla F_i(\mathbf{w}^{U,S}_t) = \nabla F^U(\mathbf{w}^{U,S}_t). \end{aligned}$$

### A.1.3 Classical results for random variables

To analyze DeltaGrad with SGD, Bernstein's inequality (Oliveira, 2009; Tropp, 2012, 2016, e.g.,) is necessary. Both its scalar version and matrix version are stated below.

**Lemma 4** (Bernstein's inequality for scalars). *Consider a list of independent random variables,  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_k$  satisfying  $\mathbb{E}(\mathbf{S}_i) = \mathbf{0}$  and  $|\mathbf{S}_i| \leq J$ , and their sum  $\mathbf{Z} = \sum_{i=1}^k \mathbf{S}_i$ . Then the following inequality holds:*

$$Pr(\|\mathbf{Z}\| \geq x) \leq \exp\left(\frac{-x^2}{\sum_{i=1}^k \mathbb{E}(\mathbf{S}_i^2) + \frac{Jx}{3}}\right), \forall x \geq 0.$$

**Lemma 5** (Bernstein's inequality for matrices). *Consider a list of independent  $d_1 \times d_2$  random matrices,  $\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_k$  satisfying  $E(\mathbf{S}_i) = \mathbf{0}$  and  $\|\mathbf{S}_i\| \leq J$ , and their sum  $\mathbf{Z} = \sum_{i=1}^k \mathbf{S}_i$ . Define the deterministic "variance surrogate":*

$$V(\mathbf{Z}) = \max\left(\left\|\sum_{i=1}^k \mathbb{E}(\mathbf{S}_i \mathbf{S}_i^*)\right\|, \left\|\sum_{i=1}^k \mathbb{E}(\mathbf{S}_i^* \mathbf{S}_i)\right\|\right). \quad (8)$$

Then the following inequalities hold:

$$\Pr(\|\mathbf{Z}\| \geq x) \leq (d_1 + d_2) \exp\left(\frac{-x^2}{V(\mathbf{Z}) + \frac{Jx}{3}}\right), \forall x \geq 0, \quad (9)$$

$$\mathbb{E}(\|\mathbf{Z}\|) \leq \sqrt{2V(\mathbf{Z}) \log(d_1 + d_2)} + \frac{1}{3}J \log(d_1 + d_2). \quad (10)$$

## A.2 Results for deterministic gradient descent

The main result for DeltaGrad with GD is Theorem 5, proved in Section A.2.8.

### A.2.1 Quasi-Newton

By following equations 1.2 and 1.3 in Byrd et al. (1994), the Quasi-Hessian update can be written as:

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \frac{\mathbf{B}_t \Delta w_t \Delta w_t^T \mathbf{B}_t}{\Delta w_t^T \mathbf{B}_t \Delta w_t} + \frac{\Delta g_t \Delta g_t^T}{\Delta g_t^T \Delta w_t}. \quad (11)$$

We have used the indices  $k$  to index the Quasi-Hessians  $\mathbf{B}_{j_k}$ . This allows us to see that they correspond to the appropriate parameter gap  $\Delta w_{j_k}$  and gradient gap  $\Delta g_{j_k}$ . The indices  $j_k$  depend on the iteration number  $t$  in the main algorithm, and they are updated by removing the “oldest” entry, and adding  $T_0$  at every period.

DeltaGrad uses equation (11) on the prior updates:

$$\mathbf{B}_{j_{k+1}} = \mathbf{B}_{j_k} - \frac{\mathbf{B}_{j_k} \Delta w_{j_k} \Delta w_{j_k}^T \mathbf{B}_{j_k}}{\Delta w_{j_k}^T \mathbf{B}_{j_k} \Delta w_{j_k}} + \frac{\Delta g_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}, \quad (12)$$

where the initialized matrix  $\mathbf{B}_{j_0}$  is  $\mathbf{B}_{j_0} = \Delta g_{i_0}^T \Delta w_{j_0} / [\Delta w_{i_0}^T \Delta w_{j_0}] \mathbf{I}$ .

We use formulas 3.5 and 2.25 from Byrd et al. (1994) for the Quasi-Newton method, with the caveat that they use slightly different notation.

For the update rule of  $\mathbf{B}_{j_k}$ , i.e.:

$$\mathbf{B}_{j_{k+1}} = \mathbf{B}_{j_k} - \frac{\mathbf{B}_{j_k} \Delta w_{j_k} \Delta w_{j_k}^T \mathbf{B}_{j_k}}{\Delta w_{j_k}^T \mathbf{B}_{j_k} \Delta w_{j_k}} + \frac{\Delta g_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}. \quad (13)$$

There is an equivalent expression for the inverse of  $\mathbf{B}_{j_k}$  as below:

$$\mathbf{B}_{j_{k+1}}^{-1} = \left( \mathbf{I} - \frac{\Delta w_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}} \right) \mathbf{B}_{j_k}^{-1} \left( \mathbf{I} - \frac{\Delta g_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}} \right) + \frac{\Delta w_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}. \quad (14)$$

See Algorithm 1 for an overview of the L-BFGS algorithm.

---

#### Algorithm 1: Overview of L-BFGS algorithm

---

**Input** : The sequence of the model parameter differences  $\Delta W = \{\Delta w_0, \Delta w_1, \dots, \Delta w_{m-1}\}$ , the sequence of the gradient differences  $\Delta G = \{\Delta g_0, \Delta g_1, \dots, \Delta g_{m-1}\}$ , a vector  $\mathbf{v}$ , history size  $m$

**Output**: Approximate results of  $\mathbf{H}(w_m)\mathbf{v}$  at point  $w_m$ , and for some  $\mathbf{v}$ , such that  $\Delta w_i \approx w_i - w_{i-1}$  for all  $i$

- 1 Compute  $\Delta W^T \Delta W$
  - 2 Compute  $\Delta W^T \Delta G$ , get its diagonal matrix  $D$  and its lower triangular submatrix  $L$
  - 3 Compute  $\sigma = \Delta g_{m-1}^T \Delta w_{m-1} / (\Delta w_{m-1}^T \Delta w_{m-1})$
  - 4 Compute the Cholesky factorization for  $\sigma \Delta W^T \Delta W + LDL^T$  to get  $JJ^T$
  - 5 Compute  $p = \begin{bmatrix} -D^{\frac{1}{2}} & D^{-\frac{1}{2}}L^T \\ \mathbf{0} & J^T \end{bmatrix}^{-1} \begin{bmatrix} D^{\frac{1}{2}} & \mathbf{0} \\ D^{-\frac{1}{2}}L^T & J^T \end{bmatrix}^{-1} \begin{bmatrix} \Delta G^T \mathbf{v} \\ \sigma \Delta W^T \mathbf{v} \end{bmatrix}$
  - 6 **return**  $\sigma \mathbf{v} - [\Delta G \quad \sigma \Delta W] p$
-

### A.2.2 Proof that Quasi-Hessians are well-conditioned

We show that the Quasi-Hessian matrices computed by L-BFGS are well-conditioned.

**Lemma 6** (Bounds on Quasi-Hessians). *The Quasi-Hessian matrices  $\mathbf{B}_{j_k}$  are well-conditioned. There exist two positive constants  $K_1$  and  $K_2$  (depending on the problem parameters  $\mu, L$ , etc) such that for any  $t$ , any vector  $\mathbf{z}$ , and all  $k \in \{0, 1, \dots, m\}$ , the following inequality holds:*

$$K_1 \|\mathbf{z}\|^2 \leq \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} \leq K_2 \|\mathbf{z}\|^2.$$

*Proof.* We start with the lower bound. Based on equation (14),  $\|\mathbf{B}_{j_k}^{-1}\|$  can be bounded by:

$$\|\mathbf{B}_{j_{k+1}}^{-1}\| \leq \|\mathbf{I} - \frac{\Delta w_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\| \cdot \|\mathbf{B}_{j_k}^{-1}\| \cdot \|\mathbf{I} - \frac{\Delta g_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\| + \|\frac{\Delta w_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\|. \quad (15)$$

in which by using the mean value theorem,  $\|\mathbf{I} - \frac{\Delta w_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\|$  can be bounded as:

$$\begin{aligned} \|\mathbf{I} - \frac{\Delta w_{j_k} \Delta g_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\| &\leq 1 + \frac{\|\Delta w_{j_k} \Delta g_{j_k}^T\|}{\Delta g_{j_k}^T \Delta w_{j_k}} \\ &= 1 + \frac{\|\Delta w_{j_k} (\mathbf{H}_{j_k} \Delta w_{j_k})^T\|}{\Delta w_{j_k}^T \mathbf{H}_{j_k} \Delta w_{j_k}} \leq 1 + \frac{\|\Delta w_{j_k}\| \|\mathbf{H}_{j_k}\| \|\Delta w_{j_k}\|}{\mu \|\Delta w_{j_k}\|^2} \leq 1 + \frac{L}{\mu}. \end{aligned} \quad (16)$$

In addition,  $\|\frac{\Delta w_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\|$  can be bounded as:

$$\|\frac{\Delta w_{j_k} \Delta w_{j_k}^T}{\Delta g_{j_k}^T \Delta w_{j_k}}\| = \|\frac{\Delta w_{j_k} \Delta w_{j_k}^T}{\Delta w_{j_k}^T \mathbf{H}_{j_k} \Delta w_{j_k}}\| \leq \|\frac{\Delta w_{j_k}^T \Delta w_{j_k}}{\mu \Delta w_{j_k}^T \Delta w_{j_k}}\| = \frac{1}{\mu}. \quad (17)$$

So by combining Equation (16) and Equation (17), Equation (15) can be bounded by:

$$\begin{aligned} \|\mathbf{B}_{j_{k+1}}^{-1}\| &\leq (1 + \frac{L}{\mu})^2 \|\mathbf{B}_{j_k}^{-1}\| + \frac{1}{\mu} \leq (1 + \frac{L}{\mu})^{2k} \|\mathbf{B}_{j_0}^{-1}\| + \frac{1 - (1 + \frac{L}{\mu})^{2k}}{1 - (1 + \frac{L}{\mu})^2} \frac{1}{\mu} \\ &= (1 + \frac{L}{\mu})^{2k} \frac{L}{\mu} + \frac{1 - (1 + \frac{L}{\mu})^{2k}}{1 - (1 + \frac{L}{\mu})^2} \frac{1}{\mu}. \end{aligned}$$

which thus implies that  $\|\mathbf{B}_{j_k}\| \geq K_1 := \frac{1}{(1 + \frac{L}{\mu})^{2k} \frac{L}{\mu} + \frac{1 - (1 + \frac{L}{\mu})^{2k}}{1 - (1 + \frac{L}{\mu})^2} \frac{1}{\mu}}$  where  $0 \leq k \leq m$ . Recall that  $m$  is small,

(set as  $m = 2$  in the experiments). So the lower bound will not approach zero.

Then based on Equation (11), we derive an upper bound for  $\|\mathbf{B}_{j_k}\|$  as follows:

$$\begin{aligned} \mathbf{z}^T \mathbf{B}_{j_{k+1}} \mathbf{z} &= \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} - \frac{\mathbf{z}^T \mathbf{B}_{j_k} \Delta w_{j_k} \Delta w_{j_k}^T \mathbf{B}_{j_k} \mathbf{z}}{\Delta w_{j_k}^T \mathbf{B}_{j_k} \Delta w_{j_k}} + \frac{\mathbf{z}^T \Delta g_{j_k} \Delta g_{j_k}^T \mathbf{z}}{\Delta g_{j_k}^T \Delta w_{j_k}} \\ &\leq \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} + \frac{\mathbf{z}^T \Delta g_{j_k} \Delta g_{j_k}^T \mathbf{z}}{\Delta g_{j_k}^T \Delta w_{j_k}} = \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} + \frac{\mathbf{z}^T \mathbf{H}_{j_k} \Delta w_{j_k} \Delta w_{j_k}^T \mathbf{H}_{j_k} \mathbf{z}}{\Delta w_{j_k}^T \mathbf{H}_{j_k} \Delta w_{j_k}} \\ &\leq \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} + \frac{\mathbf{z}^T \mathbf{H}_{j_k} \mathbf{z} \Delta w_{j_k}^T \mathbf{H}_{j_k} \Delta w_{j_k}}{\Delta w_{j_k}^T \mathbf{H}_{j_k} \Delta w_{j_k}} = \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} + \mathbf{z}^T \mathbf{H}_{j_k} \mathbf{z} \\ &\leq \mathbf{z}^T \mathbf{B}_{j_k} \mathbf{z} + L \|\mathbf{z}\|^2. \end{aligned}$$

The first inequality uses the fact that  $\mathbf{z}^T \mathbf{B}_{j_k} \Delta w_{j_k} \Delta w_{j_k}^T \mathbf{B}_{j_k} \mathbf{z} = (\mathbf{z}^T \mathbf{B}_{j_k} \Delta w_{j_k})^2 \geq 0$  and  $\Delta w_{j_k}^T \mathbf{B}_{j_k} \Delta w_{j_k} \geq 0$ , due to the positive definiteness of  $\mathbf{B}_{j_k}$ . The second inequality uses the Cauchy-Schwarz inequality for the Quasi-Hessian, i.e.:

$$(\mathbf{a}^T \mathbf{H}_{j_k} \mathbf{b})^2 \leq (\mathbf{a}^T \mathbf{H}_{j_k} \mathbf{a}) (\mathbf{b}^T \mathbf{H}_{j_k} \mathbf{b}).$$

By applying the formula above recursively, we get  $\mathbf{z}^T \mathbf{B}_{j_{k+1}} \mathbf{z} \leq (k+1)L \|\mathbf{z}\|^2$  where  $0 \leq k \leq m$ . Again, as  $m$  is bounded, so we have  $(k+1)L \leq K_2 := (m+1)L$ . This finishes the proof.  $\square$

### A.2.3 Proof preliminaries

First of all, we provide the bound on  $\delta_t$ , which is defined as:

**Lemma 7** (Upper bound on  $\delta_t$ ). *By defining*

$$\delta_t = -\frac{\eta}{n-r} \left( \frac{r}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right),$$

we then have  $\|\delta_t\| \leq 2c_2 \frac{r\eta}{n}$ .

*Proof.* Based on the definition of  $\delta_t$ , we can rearrange it a little bit as:

$$\begin{aligned} \|\delta_t\| &= \left\| -\frac{\eta r}{n(n-r)} \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) + \frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right\| \\ &= \left\| -\frac{\eta r}{n(n-r)} \left[ \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right] + \left( \frac{\eta}{n-r} - \frac{\eta r}{n(n-r)} \right) \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right\| \\ &= \left\| -\frac{\eta r}{n(n-r)} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t) + \frac{\eta}{n} \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right\|. \end{aligned}$$

Then by using the triangle inequality and Assumption 3 (bounded gradients), the formula above can be bounded as:

$$\leq \frac{\eta r}{n(n-r)} \sum_{i \notin R} \|\nabla F_i(\mathbf{w}^U_t)\| + \frac{\eta}{n} \sum_{i \in R} \|\nabla F_i(\mathbf{w}^U_t)\| \leq \frac{\eta r}{n} c_2 + \frac{\eta r}{n} c_2 = \frac{2\eta r}{n} c_2$$

□

Notice that Algorithm 1 requires  $2m$  vectors as the input, i.e.  $[\Delta w_{j_0}, \Delta w_{j_1}, \dots, \Delta w_{j_{m-1}}]$  and  $[\Delta g_{j_0}, \Delta g_{j_1}, \dots, \Delta g_{j_{m-1}}]$  to approximate the product of the Hessian matrix  $\mathbf{H}(w_t)$  and the input vector  $\Delta w_t$  at the  $t$ th iteration where  $j_{m-1} \leq t \leq j_{m-1} + T_0$ .

Note that by multiplying  $\Delta w_{j_k}$  on both sides of the Quasi-Hessian update Equation (12), we have the classical *secant equation* that characterizes Quasi-Newton methods as below:

$$\mathbf{B}_{j_{k+1}} \Delta w_{j_k} = \Delta g_{j_k}. \quad (18)$$

Then we give a bound on the quantity  $\|\Delta g_{j_k} - \mathbf{B}_{j_q} \Delta w_{j_k}\|$  where the intermediate index  $q$  is in between the "correct" index  $k+1$  and the final index  $m$ , so  $m \geq q \geq k+1$ . This characterizes the error by using a different Quasi-Hessian at some iteration. Its proof borrows ideas from Conn et al. (1991). Unlike Conn et al. (1991), our proof relies on a preliminary estimate on the bound on  $\|\mathbf{w}_t - \mathbf{w}^I_t\|$ , which is at the level of  $O(\frac{r}{n})$ . The proof of the bound will be presented later.

**Theorem 1.** *Suppose that the preliminary estimate is:  $\|\mathbf{w}_{j_k} - \mathbf{w}^I_{j_k}\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$ , where  $k = 1, 2, \dots, m$  and  $M_1 = \frac{2c_2}{\mu}$ . Let  $e = \frac{L(L+1)+K_2L}{\mu K_1}$ , for the upper and lower bounds  $K_1, K_2$  on the eigenvalues of the quasi-Hessian from Lemma 6, for the upper bounds  $c_2$  on the gradient from Assumption 3 and for the Lipschitz constant  $c_0$  of the Hessian. For  $1 \leq k+1 \leq q \leq m$ , we have:*

$$\|\mathbf{H}_{j_k} - \mathbf{H}_{j_q}\| \leq c_0 d_{j_k, j_q} + c_0 \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$$

and

$$\|\Delta g_{j_k} - \mathbf{B}_{j_q} \Delta w_{j_k}\| \leq [(1+e)^{q-k-1} - 1] \cdot c_0 (d_{j_k, j_q} + \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}) \cdot s_{j_1, j_m},$$

where  $s_{j_1, j_m} = \max(\|\Delta w_a\|)_{a=j_1, j_2, \dots, j_m}$  and  $d$  is defined as the maximum gap between the steps of the algorithm over the iterations from  $j_k$  to  $j_q$ :

$$d_{j_k, j_q} = \max(\|\mathbf{w}_a - \mathbf{w}_b\|)_{j_k \leq a \leq b \leq j_q}. \quad (19)$$

*Proof.* Let  $v_q = \Delta g_{j_k} - \mathbf{B}_{j_{q+1}} \Delta w_{j_k}$ ,  $b_q = \|v_q\|$  and  $f = c_0(d_{j_1, j_m + T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}) s_{j_1, j_m}$ .

Let us bound the difference between the averaged Hessians  $\|\mathbf{H}_{j_k} - \mathbf{H}_{j_q}\|$ , where  $1 \leq k < q \leq m$ , using their definition, as well as using Assumption 4 on the Lipschitzness of the Hessian:

$$\begin{aligned}
& \|\mathbf{H}_{j_k} - \mathbf{H}_{j_q}\| \\
&= \left\| \int_0^1 [\mathbf{H}(\mathbf{w}_{j_k} + x(\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}))] dx - \int_0^1 [\mathbf{H}(\mathbf{w}_{j_q} + x(\mathbf{w}^I_{j_q} - \mathbf{w}_{j_q}))] dx \right\| \\
&= \left\| \int_0^1 [\mathbf{H}(\mathbf{w}_{j_k} + x(\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}))] - \mathbf{H}(\mathbf{w}_{j_q} + x(\mathbf{w}^I_{j_q} - \mathbf{w}_{j_q}))] dx \right\| \\
&\leq c_0 \int_0^1 \|\mathbf{w}_{j_k} + x(\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}) - [\mathbf{w}_{j_q} + x(\mathbf{w}^I_{j_q} - \mathbf{w}_{j_q})]\| dx \\
&\leq c_0 \|\mathbf{w}_{j_k} - \mathbf{w}_{j_q}\| + \frac{c_0}{2} \|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k} - (\mathbf{w}^I_{j_q} - \mathbf{w}_{j_q})\| \\
&\leq c_0 \|\mathbf{w}_{j_k} - \mathbf{w}_{j_q}\| + \frac{c_0}{2} \|\mathbf{w}_{j_q} - \mathbf{w}^I_{j_q}\| + \frac{c_0}{2} \|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}\| \\
&\leq c_0 d_{j_k, j_q} + \frac{c_0}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n} \leq c_0 d_{j_1, j_m + T_0 - 1} + \frac{c_0}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.
\end{aligned} \tag{20}$$

On the last line, we used the definition of  $d_{j_k, j_q}$ , and the assumption on the boundedness of  $\|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}\|$ .

Then, when  $q = k$ , according to Equation (18), the secant equation  $\Delta g_{j_k} = \mathbf{B}_{j_{k+1}} \Delta w_{j_k}$  holds. So  $\|\Delta g_{j_k} - \mathbf{B}_{j_{k+1}} \Delta w_{j_k}\| = 0$ , which proves the claim when  $q = k$ . So  $v_q = b_q = 0$ .

Next, let  $u_q = \Delta g_{j_q} - \mathbf{B}_{j_q} \Delta w_{j_q}$ . This quantity is closely related to  $v_{q-1} = \Delta g_{j_k} - \mathbf{B}_{j_q} \Delta w_{j_k}$ , and the difference is that in  $u_q$ , the  $\Delta g, \Delta w$  terms are defined at  $q$ , as opposed to the base one at  $k$ . Then  $|u_q^T \Delta w_{j_k}|$ , where  $q > k$ , can be bounded as:

$$\begin{aligned}
& |u_q^T \Delta w_{j_k}| \\
&= |\Delta g_{j_q}^T \Delta w_{j_k} - \Delta g_{j_q}^T \Delta w_{j_q} + \Delta g_{j_q}^T \Delta w_{j_q} - \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k}| \\
&\leq |\Delta g_{j_q}^T \Delta w_{j_k} - \Delta g_{j_q}^T \Delta w_{j_q}| + |\Delta w_{j_q}^T v_{q-1}| \\
&\leq |\Delta g_{j_q}^T \Delta w_{j_k} - \Delta g_{j_q}^T \Delta w_{j_q}| + \|\Delta w_{j_q}\| \cdot b_{q-1} \\
&= |\Delta w_{j_q}^T \mathbf{H}_{j_q} \Delta w_{j_k} - \Delta w_{j_q}^T \mathbf{H}_{j_q} \Delta w_{j_q}| + \|\Delta w_{j_q}\| \cdot b_{q-1} \\
&= |\Delta w_{j_q}^T (\mathbf{H}_{j_q} - \mathbf{H}_{j_k}) \Delta w_{j_k}| + \|\Delta w_{j_q}\| \cdot b_{q-1} \\
&\leq \|\Delta w_{j_q}\| \cdot \|\mathbf{H}_{j_q} - \mathbf{H}_{j_k}\| \cdot \|\Delta w_{j_k}\| + \|\Delta w_{j_q}\| \cdot b_{q-1} \\
&\leq (f + b_{q-1}) \|\Delta w_{j_q}\|,
\end{aligned} \tag{21}$$

in which the first inequality uses the triangle inequality, the second inequality uses the Cauchy-Schwarz inequality, and the subsequent equality uses the Cauchy mean value theorem. Finally, the third inequality uses Assumption 4 and equation (20). We also use the following bounds, which hold by definition (notice that  $k, q \leq m$ ):

$$\|w_{j_k} - w_{j_q}\| \leq d_{j_k, j_q} \quad \|\Delta w_{j_q}\| \leq s_{j_1, j_m}.$$

The argument on the upper bound of  $b_q$  will proceed by induction. The claim is true for the base case  $q = k$ . Assuming that the claim is true for  $q - 1$ , we want to prove it for  $q$ , which is bounded as below:

$$b_q = \left\| \Delta g_{j_k} - \left( \mathbf{B}_{j_q} - \frac{\mathbf{B}_{j_q} \Delta w_{j_q} \Delta w_{j_q}^T \mathbf{B}_{j_q}}{\Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_q}} + \frac{\Delta g_{j_q} \Delta g_{j_q}^T}{\Delta g_{j_q}^T \Delta w_{j_q}} \right) \Delta w_{j_k} \right\|. \tag{22}$$

By using the triangle inequality, we obtain the following upper bound:

$$\leq b_{q-1} + \left\| \left( \frac{\Delta g_{j_q} \Delta g_{j_q}^T}{\Delta g_{j_q}^T \Delta w_{j_q}} - \frac{\mathbf{B}_{j_q} \Delta w_{j_q} \Delta w_{j_q}^T \mathbf{B}_{j_q}}{\Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_q}} \right) \Delta w_{j_k} \right\|.$$

Now we come to a key and nontrivial step of the argument. By bringing fractions to the common denominator in the second term, adding and subtracting  $\Delta g_{j_q} \Delta g_{j_q}^T \Delta w_{j_q}^T \Delta g_{j_q}$  and  $\Delta g_{j_q} (\mathbf{B}_{j_q} \Delta w_{j_q})^T \Delta w_{j_q}^T \Delta g_{j_q}$ , and rearranging to factor out the term  $-u_q$  in the numerator of each summand, the formula above can be rewritten as:

$$= b_{q-1} + \frac{\|[-\Delta g_{j_q} \Delta g_{j_q}^T \Delta w_{j_q}^T u_q + \Delta g_{j_q} u_q^T \Delta w_{j_q}^T \Delta g_{j_q} + u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_q}^T \Delta g_{j_q}] \Delta w_{j_k}\|}{\Delta g_{j_q}^T \Delta w_{j_q} \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_q}}.$$

Next, using the Cauchy mean value theorem, and the fact that the smallest eigenvalues of  $\mathbf{H}_{j_q}, \mathbf{B}_{j_q}$  are lower bounded by  $\mu, K_1$  respectively, the formula above is bounded as:

$$\begin{aligned} &\leq b_{q-1} + \frac{\|[-\Delta g_{j_q} \Delta g_{j_q}^T \Delta w_{j_q}^T u_q + \Delta g_{j_q} u_q^T \Delta w_{j_q}^T \Delta g_{j_q} + u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_q}^T \Delta g_{j_q}] \Delta w_{j_k}\|}{\mu K_1 \|\Delta w_{j_q}\|^4} \\ &\leq b_{q-1} + (\|\Delta g_{j_q}\|^2 \cdot \|\Delta w_{j_q}^T u_q \Delta w_{j_k}\| + \|\Delta g_{j_q}\| \cdot \|u_q^T \Delta w_{j_q}^T \Delta g_{j_q} \Delta w_{j_k}\| \\ &\quad + \|u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k} \Delta w_{j_q}^T \Delta g_{j_q}\|) / \mu K_1 \|\Delta w_{j_q}\|^4. \end{aligned}$$

Now we want to bound the last three terms one by one. First of all,  $\|\Delta g_{j_q}\|^2 \|\Delta w_{j_q}^T u_q \Delta w_{j_k}\|$  can be bounded as:

$$\begin{aligned} &\|\Delta g_{j_q}\|^2 \cdot \|\Delta w_{j_q}^T u_q \Delta w_{j_k}\| = \|\mathbf{H}_{j_q} \Delta w_{j_q}\|^2 \cdot |\Delta w_{j_q}^T u_q| \cdot \|\Delta w_{j_q}\| \\ &\leq L \|\Delta w_{j_q}\|^3 \cdot |\Delta w_{j_q}^T u_q| \leq L(f + b_{q-1}) \|\Delta w_{j_q}\|^4, \end{aligned}$$

in which the first equality uses the Cauchy mean value theorem, the subsequent inequality uses Assumption 3 and the last inequality uses equation (21), the upper bound on  $|\Delta w_{j_q}^T u_q|$ .

Then for  $\|\Delta g_{j_q}\| \cdot \|u_q^T \Delta w_{j_q}^T \Delta g_{j_q} \Delta w_{j_k}\|$ , we have a very similar argument. The only difference is that we factor out the scalar  $\Delta w_{j_q}^T \Delta g_{j_q}$ , and bound it by  $L \|\Delta w_{j_q}\|^2$ , i.e.:

$$\begin{aligned} &\|\Delta g_{j_q}\| \cdot \|u_q^T \Delta w_{j_q}^T \Delta g_{j_q} \Delta w_{j_k}\| \\ &= \|\mathbf{H}_{j_q} \Delta w_{j_q}\| \cdot |\Delta w_{j_q}^T \Delta g_{j_q}| \cdot |u_q^T \Delta w_{j_k}| \\ &\leq L^2 (f + b_{q-1}) \|\Delta w_{j_q}\|^4, \end{aligned}$$

in which the first equality uses Cauchy mean value theorem and the fact that  $\Delta w_{j_q}^T \Delta g_{j_q}$  is a scalar and the last inequality uses Assumption 3 and Equation (21).

In terms of the bound on  $\|u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k} \Delta w_{j_q}^T \Delta g_{j_q}\|$ , it is derived as:

$$\begin{aligned} &\|u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k} \Delta w_{j_q}^T \Delta g_{j_q}\| \\ &= \|u_q \Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k} \Delta w_{j_q}^T \Delta g_{j_q}\| \\ &\leq \|u_q \Delta w_{j_q}^T\| \cdot |\Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k}| \cdot \|\Delta g_{j_q}\| \\ &\leq (f + b_{q-1}) \|\Delta w_{j_q}\| \cdot |\Delta w_{j_q}^T \mathbf{B}_{j_q} \Delta w_{j_k}| \cdot \|\mathbf{H}_{j_q} \Delta w_{j_q}\| \\ &\leq (f + b_{q-1}) \|\Delta w_{j_q}\| \cdot K_2 \|\Delta w_{j_q}\|^2 \cdot L \|\Delta w_{j_q}\| \\ &= K_2 L (f + b_{q-1}) \|\Delta w_{j_q}\|^4, \end{aligned}$$

in which the first inequality uses the Cauchy Schwarz inequality, the second inequality uses equation (21) and the third inequality uses Assumption 6.

In summary, for all  $j \geq t + 1$ , Equation (22) is bounded by:

$$\begin{aligned} b_q &\leq b_{q-1} + \frac{L(L+1) + K_2 L}{\mu K_1 \|\Delta w_{j_q}\|^4} (f + b_{q-1}) \|\Delta w_{j_q}\|^4 \\ &= (1 + e) b_{q-1} + ef. \end{aligned}$$



By recursion and using the fact that  $b_k = 0$ , this can be bounded as:

$$\begin{aligned} &\leq (1+e)^{q-k} b_{k+1} + \sum_{i=0}^{q-k-1} (1+e)^i e \cdot f \\ &= \frac{(1+e)^{q-k} - 1}{e} \cdot ef = [(1+e)^{q-k} - 1]f. \end{aligned} \quad (23)$$

This proves the required claim  $b_q \leq [(1+e)^{q-k} - 1]f$  and finishes the proof.  $\square$

**Corollary 1** (Approximation accuracy of quasi-Hessian to mean Hessian). *Suppose that  $\|\mathbf{w}_{j_s} - \mathbf{w}^I_{j_s}\| \leq \frac{1}{2} \frac{r}{n} M_1 \frac{r}{n}$  and  $\|\mathbf{w}_t - \mathbf{w}^I_t\| \leq \frac{1}{2} \frac{r}{n} M_1 \frac{r}{n}$  where  $s = 1, 2, \dots, m$ . Then for  $j_m \leq t \leq j_m + T_0 - 1$ ,*

$$\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m} := Ad_{j_1, j_m + T_0 - 1} + A \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}, \quad (24)$$

where recall again that  $c_0$  is the Lipschitz constant of the Hessian,  $d_{j_1, j_m + T_0 - 1}$  is the maximal gap between the iterates of the GD algorithm on the full data from  $j_1$  to  $j_m + T_0 - 1$  (see equation (19)), which goes to zero as  $t \rightarrow \infty$ ) and  $A = \frac{c_0 \sqrt{m} [(1+e)^m - 1]}{c_1} + c_0$  in which  $e$  is a problem dependent constant defined in Theorem 1,  $c_1$  is the ‘‘strong independence’’ constant from (5).

*Proof.* Based on Theorem 1,  $b_{q-1} = \|\mathbf{H}_{j_q} \Delta w_{j_k} - \mathbf{B}_{j_q} \Delta w_{j_k}\| \leq [(1+e)^{q-k-1} - 1]f$ .

Then based on the ‘‘strong linear independence’’ in Assumption 5, the matrix  $\Delta W_{j_1, j_2, \dots, j_m} = [\frac{\Delta w_{j_1}}{s_{j_1, j_m}}, \frac{\Delta w_{j_2}}{s_{j_1, j_m}}, \dots, \frac{\Delta w_{j_m}}{s_{j_1, j_m}}]$  has its smallest singular value lower bounded by  $c_1 > 0$ . Then  $\|\mathbf{H}_{j_m} - \mathbf{B}_{j_m}\|$  can be bounded as below:

$$\begin{aligned} \|\mathbf{H}_{j_m} - \mathbf{B}_{j_m}\| &\leq \frac{1}{c_1} \|(\mathbf{H}_{j_m} - \mathbf{B}_{j_m}) \Delta W_{j_1, j_2, \dots, j_m}\| \\ &\leq \sqrt{m} [(1+e)^m - 1] \frac{c_0}{c_1} \left( d_{j_1, j_m + T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n} \right) \end{aligned} \quad (25)$$

The second inequality uses the bound  $\|M\| \leq \sqrt{m} \max_i \|m_i\|$ , where  $M$  is a matrix with the  $m$  columns  $m_i$ .

So by combining the results from equation (25), we can upper bound  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\|$  where  $j_m \leq t \leq j_m + T_0 - 1$ , i.e.:

$$\begin{aligned} \|\mathbf{H}_t - \mathbf{B}_{j_m}\| &= \|\mathbf{H}_t - \mathbf{H}_{j_m} + \mathbf{H}_{j_m} + \mathbf{B}_{j_m}\| \\ &\leq \|\mathbf{H}_t - \mathbf{H}_{j_m}\| + \|\mathbf{H}_{j_m} - \mathbf{B}_{j_m}\| \\ &\leq c_0 (d_{j_m, t} + M_1 \frac{r}{n}) + \sqrt{m} [(1+e)^m - 1] \frac{c_0}{c_1} \left( d_{j_1, j_m + T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n} \right) \\ &\leq Ad_{j_1, j_m + T_0 - 1} + A \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n} \end{aligned} \quad (26)$$

This finishes the proof.  $\square$

Note that in the upper bound on  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\|$ , there is one term  $d_{j_1, j_m + T_0 - 1}$ . So we need to do some analysis of this term:

**Lemma 8** (Contraction of the GD iterates). *Recall the definition of  $d_{j_k, j_q}$  from Theorem 1:*

$$d_{j_k, j_q} = \max (\|\mathbf{w}_a - \mathbf{w}_b\|)_{j_k \leq a \leq b \leq j_q}.$$

Then  $d_{j_k, j_q} \leq d_{j_k - z, j_q - z}$  for any positive integers  $z$  and  $d_{j_k, j_q} \leq (1 - \mu\eta)^{j_k} d_{0, j_q - j_k}$  for any  $0 \leq j_k \leq j_q$ .

*Proof.* To prove the two inequalities, we should look at  $d_{j_k, j_q}$  and  $d_{j_k-z, j_q-z}$  where  $z$  is a positive integer. For any given  $j_k \leq a \leq b \leq j_q$ , the upper bound on  $\|w_a - w_b\|$  can be derived as below:

$$\begin{aligned}
\|\mathbf{w}_a - \mathbf{w}_b\| &= \|\mathbf{w}_{a-1} - \eta \nabla F(\mathbf{w}_{a-1}) - (\mathbf{w}_{b-1} - \eta \nabla F(\mathbf{w}_{b-1}))\| \\
&= \|\mathbf{w}_{a-1} - \mathbf{w}_{b-1} - \eta(\nabla F(\mathbf{w}_{a-1}) - \nabla F(\mathbf{w}_{b-1}))\| \\
&= \|\mathbf{w}_{a-1} - \mathbf{w}_{b-1} - \\
&\quad \eta \frac{1}{n} \left( \int_0^1 \sum_{i=1}^n \mathbf{H}_i(\mathbf{w}_{a-1} + x(\mathbf{w}_{b-1} - \mathbf{w}_{a-1})) dx \right) (\mathbf{w}_{a-1} - \mathbf{w}_{b-1})\| \\
&= \left\| \left( \mathbf{I} - \frac{\eta}{n} \left( \int_0^1 \sum_{i=1}^n \mathbf{H}_i(\mathbf{w}_{a-1} + x(\mathbf{w}_{b-1} - \mathbf{w}_{a-1})) dx \right) \right) (\mathbf{w}_{a-1} - \mathbf{w}_{b-1}) \right\|.
\end{aligned}$$

The derivation above uses the update rule of gradient descent and Cauchy mean-value theorem. Then according to Cauchy Schwarz inequality and strong convexity, it can be further bounded as  $\|\mathbf{w}_a - \mathbf{w}_b\| \leq (1 - \eta\mu)\|\mathbf{w}_{a-1} - \mathbf{w}_{b-1}\|$ .

This can be used iteratively, which ends up with the following inequality:

$$\|\mathbf{w}_a - \mathbf{w}_b\| \leq (1 - \eta\mu)^z \|\mathbf{w}_{a-z} - \mathbf{w}_{b-z}\| \quad (27)$$

which indicates that  $d_{j_k, j_q} \leq (1 - \eta\mu)^z d_{j_k-z, j_q-z}$  and thus  $d_{j_k, j_q} \leq d_{j_k-z, j_q-z}$ . So by replacing  $z$  with  $j_k$ , we will have:  $d_{j_k, j_q} \leq (1 - \mu\eta)^{j_k} d_{0, j_q-j_k}$ .  $\square$

#### A.2.4 Main recursions

We bound the difference between  $\mathbf{w}_t^I$  and  $\mathbf{w}_t^U$ . The proofs of the theorems stated below are in the following sections.

Our proof starts out with the usual approach of trying to show a contraction for the gradient updates, see e.g., Bottou et al. (2018). First we bound  $\|\mathbf{w}_t - \mathbf{w}_t^U\|$ , i.e.:

**Theorem 2** (Bound between iterates on full and the leave- $r$ -out dataset).  $\|\mathbf{w}_t - \mathbf{w}_t^U\| \leq M_1 \frac{r}{n}$  where  $M_1 = \frac{2}{\mu} c_2$  is some positive constant that does not depend on  $t$ .

To show that the preliminary estimate on the bound on  $\|\mathbf{w}_t^I - \mathbf{w}_t\|$  used in Theorem 1 and Corollary 1 holds, the proof is provided as below:

**Theorem 3** (Bound between iterates on full data and incrementally updated ones). *Consider an iteration  $t$  indexed with  $j_m$  for which  $j_m \leq t < j_m + T_0 - 1$ , and suppose that we are at the  $x$ -th iteration of full gradient updates, so  $j_1 = j_0 + xT_0$ ,  $j_m = j_0 + (m-1+x)T_0$ . Suppose that we have the bounds  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m} = A d_{j_1, j_m + T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{r}{n}} A M_1 \frac{r}{n}$  (where we recalled the definition of  $\xi$ ) and  $\xi_{j_1, j_m} \leq \frac{\mu}{2}$  for all iterations  $x$ . Then*

$$\|\mathbf{w}_{t+1}^I - \mathbf{w}_{t+1}\| \leq \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_0, j_0+(m-1)T_0}} \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.$$

Recall that  $c_0$  is the Lipschitz constant of the Hessian,  $M_1$  and  $A$  are defined in Theorem 2 and Corollary 1 respectively, which do not depend on  $t$ ,

For this theorem, note that this inequality depends on the condition  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$  while in Theorem 1, to prove  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$ , we need to use the inequality in Theorem 3, i.e.  $\|\mathbf{w}_{t+1}^I - \mathbf{w}_{t+1}\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$ . In what follows, we will show that both inequalities hold for all the iterations  $t$  without relying on other conditions.

We can select hyper-parameters  $T_0, j_0$  such that

$$A(1 - \eta\mu)^{j_0-m+1} d_{0, (m-1)T_0} + \frac{1}{\frac{1}{2} - \frac{r}{n}} A M_1 \frac{r}{n} < \min\left(\frac{\mu}{2}, (1 - \frac{r}{n})\mu - \frac{c_0 M_1 r(n-r)}{2n^2}\right),$$

e.g. when  $m = 2$  and  $T_0 = 5$ , which is what we used in our experiments. It is enough that

$$j_0 > \max\left(\frac{\log\left(\frac{1}{Ad_{0,5}}\left[\frac{\mu}{2} - \frac{1}{\frac{1}{2}-\frac{r}{n}}AM_1\frac{r}{n}\right]\right)}{\log(1-\eta\mu)}, \frac{\log\left(\frac{1}{Ad_{0,5}}\left[(1-\frac{r}{n})\mu - \frac{1}{\frac{1}{2}-\frac{r}{n}}AM_1\frac{r}{n}\right]\right)}{\log(1-\eta\mu)}\right) + m - 1.$$

This holds for small enough  $r/n$ :

$$j_0 > \frac{\log\left(\frac{1}{Ad_{0,5}}\left[\frac{\mu}{2} - \frac{1}{\frac{1}{2}-\frac{r}{n}}AM_1\frac{r}{n}\right]\right)}{\log(1-\eta\mu)} + m - 1$$

Then the following two theorems hold.

**Theorem 4** (Bound between iterates on full data and incrementally updated ones (all iterations)). *For any  $j_m < t < j_m + T_0 - 1$ ,  $\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}$  and  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$ .*

Then we have the following bound for  $\|\mathbf{w}^U_t - \mathbf{w}^I_t\|$ , which is our main result.

**Theorem 5** (Convergence rate of DeltaGrad). *For all iterations  $t$ , the result  $\mathbf{w}^I_t$  of DeltaGrad, Algorithm 1, approximates the correct iteration values  $\mathbf{w}^U_t$  at the rate*

$$\|\mathbf{w}^U_t - \mathbf{w}^I_t\| = o\left(\frac{r}{n}\right).$$

So  $\|\mathbf{w}^U_t - \mathbf{w}^I_t\|$  is of a lower order than  $\frac{r}{n}$ .

This is proved in Section A.2.8.

### A.2.5 Proof of Theorem 2

*Proof.* By subtracting the GD update from equation (1), we have:

$$\begin{aligned} \mathbf{w}^U_{t+1} - \mathbf{w}_{t+1} &= \mathbf{w}^U_t - \mathbf{w}_t \\ &- \eta \left( \frac{1}{n-r} \left( \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right) - \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}_t) \right) \end{aligned} \quad (28)$$

in which the right-hand side can be rewritten as:

$$\begin{aligned} &\mathbf{w}^U_t - \mathbf{w}_t - \eta (\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t)) \\ &- \eta \left( \frac{1}{n-r} \left( \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right) - \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) \right) \\ &= \mathbf{w}^U_t - \mathbf{w}_t - \eta (\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t)) \\ &- \frac{\eta}{n-r} \left( \frac{r}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right) \\ &= \mathbf{w}^U_t - \mathbf{w}_t - \eta (\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t)) + \delta_t. \end{aligned}$$

Then by applying Cauchy mean value theorem, the triangle inequality, Cauchy schwarz inequality and Lemma 7 respectively, we have:

$$\begin{aligned} &\|\mathbf{w}_{t+1} - \mathbf{w}^U_{t+1}\| \\ &\leq \|\mathbf{w}_t - \mathbf{w}^U_t - \eta \left( \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx \right) (\mathbf{w}_t - \mathbf{w}^U_t)\| + \|\delta_t\| \\ &\leq \|\mathbf{I} - \eta \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx\| \|\mathbf{w}_t - \mathbf{w}^U_t\| + \frac{2c_2 r \eta}{n} \end{aligned}$$

Then by applying the triangle inequality over integrals and Lemma 2, the formula can be further bounded as:

$$\begin{aligned} &\leq \left\| \int_0^1 (\mathbf{I} - \eta \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}_t^U - \mathbf{w}_t))) dx \right\| \|\mathbf{w}_t - \mathbf{w}_t^U\| + \frac{2c_2 r \eta}{n} \\ &\leq (1 - \eta \mu) \|\mathbf{w}_t - \mathbf{w}_t^U\| + \frac{2c_2 r \eta}{n} \end{aligned}$$

Then by applying this formula iteratively, we get:

$$\|\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^U\| \leq \frac{1}{\eta \mu} \frac{2c_2 r \eta}{n} = \frac{2c_2}{\mu} \frac{r}{n} := M_1 \frac{r}{n}$$

□

### A.2.6 Proof of Theorem 3

*Proof.* The updates for the iterations  $j_m \leq t \leq j_m + T_0 - 1$  follow the Quasi-Hessian update. We proceed in a similar way as before, by expanding the recursion as below:

$$\begin{aligned} &\|\mathbf{w}_{t+1}^I - \mathbf{w}_{t+1}\| \\ &= \|\mathbf{w}_t^I - (\mathbf{w}_t - \eta \nabla F(\mathbf{w}_t)) \\ &\quad - \frac{\eta}{n-r} (n [\mathbf{B}_{j_m}(\mathbf{w}_t^I - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] - \sum_{i \in R} \nabla F_i(\mathbf{w}_t^I))\| \\ &= \|(\mathbf{I} - \eta \frac{n}{n-r} \mathbf{B}_{j_m})(\mathbf{w}_t^I - \mathbf{w}_t) - \frac{r\eta}{n-r} \nabla F(\mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}_t^I)\| \end{aligned} \quad (29)$$

By rearranging the formula above and using the triangle inequality, we get:

$$\begin{aligned} &= \|(\mathbf{I} - \eta \frac{n}{n-r} \mathbf{B}_{j_m})(\mathbf{w}_t^I - \mathbf{w}_t) - \frac{r\eta}{n-r} \nabla F(\mathbf{w}_t) \\ &\quad + \frac{\eta}{n-r} \sum_{i \in R} (\mathbf{H}_{t,i} \times (\mathbf{w}_t^I - \mathbf{w}_t) + \nabla F_i(\mathbf{w}_t))\| \\ &\leq \|(\mathbf{I} - \eta \frac{n}{n-r} \mathbf{B}_{j_m})(\mathbf{w}_t^I - \mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \in R} \mathbf{H}_{t,i} \times (\mathbf{w}_t^I - \mathbf{w}_t)\| \\ &\quad + \|\frac{r\eta}{n-r} \nabla F(\mathbf{w}_t)\| + \|\frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}_t)\| \end{aligned} \quad (30)$$

in which we use  $\mathbf{H}_{t,i}$  to denote  $\int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}_t^I - \mathbf{w}_t)) dx$  (recall that  $\mathbf{H}_i$  represents the Hessian matrix evaluated at the  $i_{th}$  sample). Then the terms in the first absolute value are rewritten as:

$$\begin{aligned} &[\mathbf{I} - \eta \frac{n}{n-r} (\mathbf{B}_{j_m} - \mathbf{H}_t + \mathbf{H}_t)] (\mathbf{w}_t^I - \mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \in R} \mathbf{H}_{t,i} \times (\mathbf{w}_t^I - \mathbf{w}_t) \\ &= [\mathbf{I} - \eta \frac{n}{n-r} (\mathbf{B}_{j_m} - \mathbf{H}_t)] (\mathbf{w}_t^I - \mathbf{w}_t) - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \times (\mathbf{w}_t^I - \mathbf{w}_t) \end{aligned}$$

which uses the fact that  $\mathbf{H}_t = \sum_{i=1}^n \mathbf{H}_{t,i} = \sum_{i \notin R} \mathbf{H}_{t,i} + \sum_{i \in R} \mathbf{H}_{t,i}$ . Then Formula (30) can be further bounded as:

$$\begin{aligned} &\leq \|[\mathbf{I} - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i}] (\mathbf{w}_t^I - \mathbf{w}_t)\| + \frac{n\eta}{n-r} \|(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}_t^I - \mathbf{w}_t)\| \\ &\quad + \|\frac{r\eta}{n-r} \nabla F(\mathbf{w}_t)\| + \|\frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}_t)\| \\ &\leq (1 - \eta \mu + \eta \frac{n}{n-r} \xi_{j_1, j_m}) \|\mathbf{w}_t^I - \mathbf{w}_t\| + \frac{r\eta c_2}{n-r} + \frac{\eta r c_2}{n-r} \end{aligned} \quad (31)$$

Then according to Lemma 8,  $d_{j_1, j_m + T_0 - 1} = d_{j_0 + xT_0, j_0 + (x+m)T_0 - 1}$  decreases with increasing  $x$ , and thus  $\xi_{j_1, j_m} = Ad_{j_1, j_m + T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{r}{n}} AM_1 \frac{r}{n}$  is also decreasing with increasing  $x$ . So the formula above can be further bounded as:

$$\leq (1 - \eta\mu + \eta \frac{n}{n-r} \xi_{j_0, j_0 + (m-1)T_0}) \|\mathbf{w}^I_t - \mathbf{w}_t\| + \frac{2r\eta c_2}{n-r}$$

This shows a recurrent inequality for  $\|\mathbf{w}^I_t - \mathbf{w}_t\|$ . Next, notice that the conditions for deriving the above inequality hold for all  $j_m \leq t \leq j_m + T_0 - 1$ .

Then, when we reach  $t = j_m$ , we have an iteration where the gradient is computed exactly. For these iterations we have  $\mathbf{w}^I_{t+1} = \mathbf{w}^I_t - \frac{\eta}{n-r} \sum_{i \notin R} \nabla F(\mathbf{w}^I_t)$  as well as  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla F(\mathbf{w}_t)$ . Using the same argument as in the bound for  $\mathbf{w}_t - \mathbf{w}^U_t$  we can get:

$$\|\mathbf{w}_{t+1} - \mathbf{w}^I_{t+1}\| \leq [1 - \eta\mu] \|\mathbf{w}_t - \mathbf{w}^I_t\| + \frac{2c_2 r \eta}{n}.$$

Therefore, we effectively have  $\xi = 0$  for these iterations. We then continue with  $t \leftarrow t - 1$ , and use the appropriate bound among the two derived above. This recursive process works until we reach  $t = 1$ .

As long as  $\xi_{j_0, j_0 + (m-1)T_0} \leq \frac{\mu}{2}$ ,  $-\eta\mu + \eta \frac{n}{n-r} \xi_{j_0, j_0 + (m-1)T_0} < -\eta\mu + \eta \frac{n}{n-r} \frac{\mu}{2} < 0$ . Then we get the following inequality:

$$\begin{aligned} \|\mathbf{w}^I_t - \mathbf{w}_t\| &\leq \frac{2 \frac{nr c_2}{n-r}}{\eta\mu - \eta \frac{n}{n-r} \xi_{j_0, j_0 + (m-1)T_0}} \\ &= \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_0, j_0 + (m-1)T_0}} \end{aligned}$$

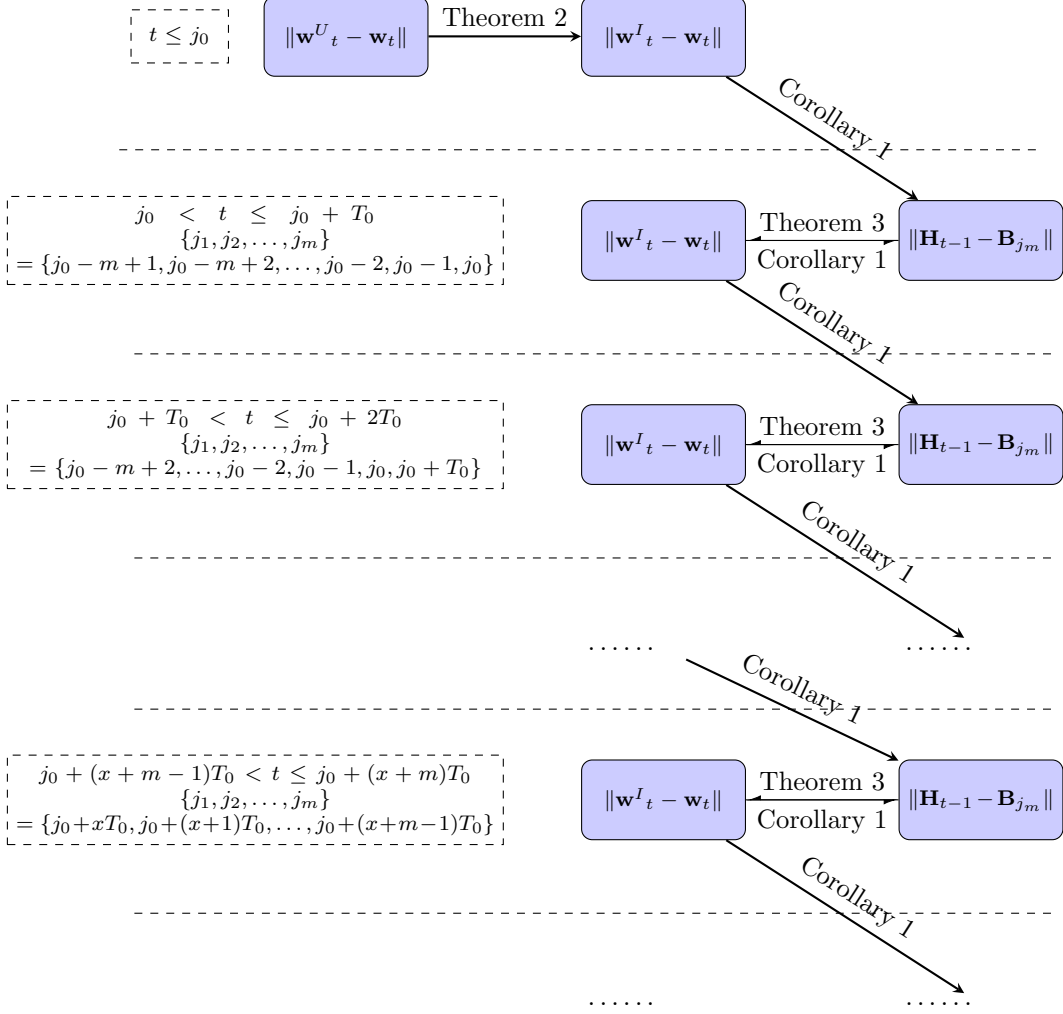
As long as  $\xi_{j_0, j_0 + (m-1)T_0} \leq \frac{\mu}{2}$ , then

$$\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_0, j_0 + (m-1)T_0}} \leq \frac{2rc_2/n}{(1-r/n)\mu - \frac{\mu}{2}} = \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.$$

The last step uses the fact that  $M_1 = \frac{2c_2}{\mu}$ . □

### A.2.7 Proof of Theorem 4

**Architecture of the proof.** To visualize the recursive proof process, we draw a picture as:



*Proof.* First of all, in terms of the bound on  $\xi_{j_1, j_m}$  which is required in Theorem 3, i.e.  $\xi_{j_1, j_m} \leq \frac{\mu}{2}$ , we do the analysis below to show that we can adjust the value of  $j_0$  and  $T_0$  such that it can hold for all  $t$ . When  $j_1 \geq j_0$ , i.e.  $j_1 = j_0 + xT_0$ , then

$$\{j_2, j_3, \dots, j_m\} = \{j_0 + (x+1)T_0, j_0 + (x+2)T_0, \dots, j_0 + (x+m-1)T_0\},$$

thus  $\xi_{j_1, j_m} = \xi_{j_0 + xT_0, j_0 + (x+m-1)T_0} = Ad_{j_0 + xT_0, j_0 + (x+m)T_0 - 1} + \frac{1}{\frac{1}{2} - \frac{x}{n}} AM_1 \frac{r}{n}$ . Here  $d_{j_0 + xT_0, j_0 + (x+m)T_0 - 1}$  decreases with  $x$ , and so does  $\xi_{j_1, j_m} = \xi_{j_0 + xT_0, j_0 + (x+m-1)T_0}$ . So the following inequality holds:

$$d_{j_0 + xT_0, j_0 + (x+m)T_0 - 1} \leq d_{j_0, j_0 + mT_0 - 1} \leq (1 - \mu\eta)^{j_0} d_{0, mT_0 - 1}.$$

When  $j_1 < j_0$ , there are only  $m$  different choices for  $\{j_1, j_2, \dots, j_m\}$ , in which the smallest  $j_1$  used for approximation is  $j_0 - m + 1$ . Then, the following inequality holds:

$$d_{j_1, j_m} \leq (1 - \eta\mu)^{j_1} d_{0, j_m - j_1} \leq (1 - \eta\mu)^{j_0 - m + 1} d_{0, j_m - j_1}.$$

For those  $j_1, j_2, \dots, j_m$ , we have  $j_m - j_1 \leq (m-1)T_0$  and thus

$$d_{j_1, j_m} \leq (1 - \eta\mu)^{j_1} d_{0, j_m - j_1} \leq (1 - \eta\mu)^{j_0 - m + 1} d_{0, j_m - j_1} \leq (1 - \eta\mu)^{j_0 - m + 1} d_{0, (m-1)T_0}.$$

So  $\xi_{j_1, j_m}$  is bounded by  $A(1 - \eta\mu)^{j_0 - m + 1} d_{0, (m-1)T_0} + \frac{1}{\frac{1}{2} - \frac{x}{n}} AM_1 \frac{r}{n}$ . To make sure  $\xi_{j_1, j_m} \leq \frac{\mu}{2}$ , we can adjust  $j_0, m, T_0$  to make  $A(1 - \eta\mu)^{j_0 - m + 1} d_{0, (m-1)T_0} + \frac{1}{\frac{1}{2} - \frac{x}{n}} AM_1 \frac{r}{n}$  smaller than  $\frac{\mu}{2}$ .

Then when  $t \leq j_0$ , the gradient is evaluated explicitly, which means that  $\mathbf{w}^U_t = \mathbf{w}^I_t$ , so the bound clearly holds, i.e., from Theorem 2, we have  $\|\mathbf{w}^U_t - \mathbf{w}_t\| \leq \frac{M_1 r}{n}$  and thus  $\|\mathbf{w}^I_t - \mathbf{w}_t\| = \|\mathbf{w}^U_t - \mathbf{w}_t\| \leq \frac{M_1 r}{n} \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$ .

When  $j_0 < t < j_0 + T_0$ , in order to compute  $\mathbf{w}^I_t$ , we need to use the history information  $\{\Delta w_{j_1}, \Delta w_{j_2}, \dots, \Delta w_{j_m}\}$ ,  $\{\Delta g_{j_1}, \Delta g_{j_2}, \dots, \Delta g_{j_m}\}$  and the corresponding quasi-Hessian matrices  $\{\mathbf{B}_{j_1}, \mathbf{B}_{j_2}, \dots, \mathbf{B}_{j_m}\}$  where  $\{j_1, j_2, \dots, j_m\} = \{j_0 - m + 1, j_0 - m + 2, \dots, j_0\}$  (we suppose  $m < j_0$ , which is a natural assumption). Since  $\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{M_1 r}{n}$  for any  $t \leq j_0$ , the conditions of Corollary 1 (used here with the  $j_1, \dots, j_m$  described above) hold up to  $j_0$ , so when  $t = j_0 + 1$ ,  $\|\mathbf{H}_{t-1} - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$  where

$$\xi_{j_1, j_m} = \xi_{j_0 - m + 1, j_0} = Ad_{j_1, j_m + T_0 - 1} + AM_1 \frac{r}{n} = Ad_{j_0 - m + 1, j_0 + T_0} + AM_1 \frac{r}{n}.$$

Plus, according to Theorem 3,  $\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_1, j_m}} = \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_0 - m + 1, j_0}}$ . When  $\xi_{j_0 - m + 1, j_0} \leq \frac{\mu}{2}$ , then

$$\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_1, j_m}} = \frac{2rc_2/n}{(1-r/n)\mu - \frac{\mu}{2}} = \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.$$

So the bound on  $\|\mathbf{w}^I_t - \mathbf{w}_t\|$  holds for all  $t \leq j_0 + 1$ . Then according to the conditions of Corollary 1, when  $t = j_0 + 2$ ,  $\|\mathbf{H}_{t-1} - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$  holds. This can proceed recursively until  $t = j_0 + T_0$ , in which the gradients are explicitly evaluated according to Theorem 3, i.e.:

$$\|\mathbf{w}^I_{j_0 + T_0} - \mathbf{w}_{j_0 + T_0}\| \leq \frac{2rc_2/n}{(1-r/n)\mu - \xi_{j_1, j_m}} = \frac{2rc_2/n}{(1-r/n)\mu - \frac{\mu}{2}} = \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.$$

Next when  $j_0 + T_0 < t < j_0 + 2T_0$ ,  $j_m$  is updated as  $j_0 + T_0$  while  $\{j_1, j_2, \dots, j_{m-1}\}$  is updated as  $\{j_0 - m + 2, j_0 - m + 3, \dots, j_0\}$  and we know that  $\|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$ . So based on Corollary 1, the following inequality holds:

$$\begin{aligned} \|\mathbf{H}_t - \mathbf{B}_{j_m}\| &\leq \xi_{j_1, j_m} = \xi_{j_0 - m + 2, j_0 + T_0} = Ad_{j_1, j_m + T_0 - 1} + AM_1 \frac{r}{n} \\ &= Ad_{j_0 - m + 2, j_0 + 2T_0} + AM_1 \frac{r}{n} \end{aligned}$$

This process can proceed recursively.

When  $j_0 + xT_0 < t < j_0 + (x+1)T_0$ , we know that:

$$\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m} = Ad_{j_1, j_m + T_0 - 1} + AM_1 \frac{r}{n}.$$

Then based on Theorem 3,  $\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$ . Then at iteration  $j_0 + (x+1)T_0$ , we update  $j_1, j_2, \dots, j_m$  as:  $j_m \leftarrow j_0 + (x+1)T_0$  and  $j_{i-1} \leftarrow j_i$  ( $i = 2, 3, \dots, m$ ) and thus

$$\|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$$

still holds for all  $k = 1, 2, \dots, m$ .

So when  $j_0 + (x+1)T_0 < t < j_0 + (x+2)T_0$ , Corollary 1 and Theorem 3 are applied alternatively. Then the following two inequalities hold for all iterations  $t$  satisfying  $j_0 + (x+1)T_0 < t < j_0 + (x+2)T_0$ :

$$\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m} = Ad_{j_1, j_m + T_0 - 1} + AM_1 \frac{r}{n},$$

$$\|\mathbf{w}^I_{j_k} - \mathbf{w}_{j_k}\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}.$$

So in the end, we know that:

$$\|\mathbf{w}^I_t - \mathbf{w}_t\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}$$

and

$$\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$$

hold for all  $t$ .

□

### A.2.8 Proof of Theorem 5

*Proof.* The proof is by induction.

When  $t \leq j_0$ , the gradient is evaluated explicitly, which means that  $\mathbf{w}^U_t = \mathbf{w}^I_t$ , so the bound clearly holds.

From iteration  $T_0$  to iteration  $t$ , the difference between  $\mathbf{w}^I_t$  and  $\mathbf{w}^U_t$  can be bounded as follows. In these equations, we use the definition of the update formula  $\mathbf{w}^I_{t+1} = \mathbf{w}^I_t - \frac{\eta}{n-r} [n(\mathbf{B}_{j_m}(\mathbf{w}^I_t - \mathbf{w}_t)) - \sum_{i \in R} \nabla F(\mathbf{w}^I_t)]$ . By rearranging terms appropriately, we get:

$$\begin{aligned} \|\mathbf{w}^I_{t+1} - \mathbf{w}^U_{t+1}\| &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{n\eta}{n-r} [\mathbf{B}_{j_m}(\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] \\ &\quad + \frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}^I_t) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)\| \end{aligned} \quad (32)$$

Then by bringing in  $\mathbf{H}_t$  into the expression above, it is rewritten as:

$$\begin{aligned} &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}^I_t - \mathbf{w}_t) + \mathbf{H}_t \times (\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] \\ &\quad + \frac{\eta}{n-r} \sum_{i \in R} (\nabla F_i(\mathbf{w}^I_t) - \nabla F_i(\mathbf{w}_t) + \nabla F_i(\mathbf{w}_t)) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)\| \end{aligned} \quad (33)$$

In the formula above, we will try to make sure there is no confusion between  $\mathbf{H}_t(\mathbf{w})$  (Hessian as a function evaluated at  $\mathbf{w}$ ) and  $\mathbf{H}_t \times (\mathbf{w})$  (Hessian times a vector). Then by applying the Cauchy mean value theorem over each individual  $\nabla F_i(\mathbf{w}^I_t) - \nabla F_i(\mathbf{w}_t)$  and by denoting the corresponding Hessian matrix as  $\mathbf{H}_{t,i}$  (note that  $\sum_{i=1}^n \mathbf{H}_{t,i} = n\mathbf{H}_t$ ), the expression becomes:

$$\begin{aligned} &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}^I_t - \mathbf{w}_t) + \mathbf{H}_t \times (\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] \\ &\quad + \frac{\eta}{n-r} \sum_{i \in R} (\mathbf{H}_{t,i} \times (\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F_i(\mathbf{w}_t)) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)\| \end{aligned}$$

Then by using the fact that  $\sum_{i \in R} \nabla F_i(\mathbf{w}_t) + \sum_{i \notin R} \nabla F_i(\mathbf{w}_t) = n\nabla F(\mathbf{w}_t)$  and  $\sum_{i \in R} \mathbf{H}_{t,i} + \sum_{i \notin R} \mathbf{H}_{t,i} = n\mathbf{H}_t$ , the expression can be rearranged as:

$$\begin{aligned} &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \times (\mathbf{w}^I_t - \mathbf{w}_t) - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}^I_t - \mathbf{w}_t)] \\ &\quad - \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)\| \end{aligned}$$

in which  $\frac{\eta}{n-r} \sum_{i \in R} \nabla F_i(\mathbf{w}_t)$  is canceled out. Then by adding and subtracting  $\mathbf{w}^U_t$  in the first part, we get:

$$\begin{aligned} &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \times (\mathbf{w}^I_t - \mathbf{w}^U_t) - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}^I_t - \mathbf{w}^U_t)] \\ &\quad - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \times (\mathbf{w}^U_t - \mathbf{w}_t) - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t)(\mathbf{w}^U_t - \mathbf{w}_t)] \\ &\quad - \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)\| \end{aligned}$$



We apply Cauchy mean value theorem over  $-\frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)$ , i.e.:

$$\begin{aligned} & -\frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}_t) + \frac{\eta}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t) \\ &= \frac{\eta}{n-r} \left[ \sum_{i \notin R} \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx \right] (\mathbf{w}^U_t - \mathbf{w}_t). \end{aligned}$$

In addition, note that  $\mathbf{H}_{t,i} = \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^I_t - \mathbf{w}_t)) dx$ . So the formula above becomes:

$$\begin{aligned} &= \|\mathbf{w}^I_t - \mathbf{w}^U_t - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \times (\mathbf{w}^I_t - \mathbf{w}^U_t) - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t) (\mathbf{w}^I_t - \mathbf{w}^U_t)] \\ & - \frac{\eta}{n-r} \sum_{i \notin R} \left( \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^I_t - \mathbf{w}_t)) dx \right) (\mathbf{w}^U_t - \mathbf{w}_t) - \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t) (\mathbf{w}^U_t - \mathbf{w}_t)] \\ & + \frac{\eta}{n-r} \sum_{i \notin R} \left( \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx \right) (\mathbf{w}^U_t - \mathbf{w}_t) \|. \end{aligned}$$

Then by applying the triangle inequality and rearranging the expression appropriately, the expression can be bounded as:

$$\begin{aligned} &\leq \left\| \left( \mathbf{I} - \frac{\eta}{n-r} \sum_{i \notin R} \mathbf{H}_{t,i} \right) (\mathbf{w}^I_t - \mathbf{w}^U_t) \right\| + \left\| \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t) (\mathbf{w}^I_t - \mathbf{w}^U_t)] \right\| \\ & + \left\| \frac{\eta}{n-r} \left[ \sum_{i \notin R} \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx - \int_0^1 \mathbf{H}_i(\mathbf{w}_t + x(\mathbf{w}^I_t - \mathbf{w}_t)) dx \right] (\mathbf{w}^U_t - \mathbf{w}_t) \right\| \\ & + \left\| \frac{n\eta}{n-r} [(\mathbf{B}_{j_m} - \mathbf{H}_t) (\mathbf{w}^U_t - \mathbf{w}_t)] \right\|, \end{aligned}$$

in which the first term is the main contraction component which always appears in the analyses of gradient descent type algorithms. The remaining terms are error terms due to the various sources of error: using a quasi-Hessian, not having a quadratic objective (implicitly assumed by the local models at each step), using the iterate  $\mathbf{w}^I$  for our update instead of the correct  $\mathbf{w}^U$ .

Then by using the following facts:

1.  $\|\mathbf{I} - \eta \mathbf{H}_{t,i}\| \leq 1 - \eta\mu$ ;
2. from Theorem 4 on the approximation accuracy of the quasi-Hessian to mean Hessian, we have the error bound  $\|\mathbf{H}_t - \mathbf{B}_{j_m}\| \leq \xi_{j_1, j_m}$ ;
3. we can bound the difference of integrated Hessians using the strategy from equation (20);
4. from Theorem 2, we have the error bound  $\|\mathbf{w}^U_t - \mathbf{w}_t\| \leq M_1 \frac{r}{n}$  (and this requires no additional assumptions),

the expression above can be bounded as follows:

$$\begin{aligned} &\leq (1 - \eta\mu + \frac{n\eta}{n-r} \xi_{j_1, j_m}) \|\mathbf{w}^I_t - \mathbf{w}^U_t\| + \frac{\eta c_0}{2} \|\mathbf{w}^U_t - \mathbf{w}^I_t\| \|\mathbf{w}^U_t - \mathbf{w}_t\| \\ & + \frac{n\eta}{n-r} \xi_{j_1, j_m} \|\mathbf{w}^U_t - \mathbf{w}_t\| \\ &\leq (1 - \eta\mu + \frac{n\eta}{n-r} \xi_{j_1, j_m} + \frac{c_0 M_1 r \eta}{2n}) \|\mathbf{w}^I_t - \mathbf{w}^U_t\| + \frac{M_1 r \eta}{n-r} \xi_{j_1, j_m} \end{aligned} \tag{34}$$

Recall from Corollary 1 that  $\xi_{j_1, j_m} = \xi_{j_0+xT_0, j_0+(m+x-1)T_0} = Ad_{j_0+xT_0, j_0+(m+x)T_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}$  decreases with the increasing  $x$ . So the formula above can be bounded as:

$$\leq (1 - \eta\mu + \frac{n\eta}{n-r}\xi_{j_0, j_0+(m-1)T_0} + \frac{c_0M_1r\eta}{2n})\|\mathbf{w}^I_t - \mathbf{w}^U_t\| + \frac{M_1r\eta}{n-r}\xi_{j_1, j_m}. \quad (35)$$

Also by plugging the formula for  $\xi$  into the formula above and using Lemma 8 (contraction of GD updates), we get:

$$\begin{aligned} &\leq (1 - \eta\mu + \frac{n\eta}{n-r}\xi_{j_0, j_0+(m-1)T_0} + \frac{c_0M_1r\eta}{2n})\|\mathbf{w}^I_t - \mathbf{w}^U_t\| \\ &+ \frac{M_1r\eta}{n-r}(Ad_{j_0+xT_0, j_0+(m+x)T_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}) \\ &\leq (1 - \eta\mu + \frac{n\eta}{n-r}\xi_{j_0, j_0+(m-1)T_0} + \frac{c_0M_1r\eta}{2n})\|\mathbf{w}^I_t - \mathbf{w}^U_t\| \\ &+ \frac{M_1r\eta}{n-r}(A(1 - \eta\mu)^{j_0+xT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}) \end{aligned} \quad (36)$$

Now, we will argue that it is possible to choose hyperparameters such that  $\xi_{j_0, j_0+(m-1)T_0} \leq (1 - \frac{r}{n})\mu - \frac{c_0M_1r(n-r)}{2n^2}$ . Then  $1 - \eta\mu + \frac{n\eta}{n-r}\xi_{j_0, j_0+(m-1)T_0} + \frac{c_0M_1r\eta}{2n}$  is a constant for all  $t$  and smaller than 1. By denoting  $\mu - \frac{n}{n-r}\xi_{j_0, j_0+(m-1)T_0} - \frac{c_0M_1r}{2n}$  as  $C$ , the formula above can be written as:

$$= (1 - \eta C)\|\mathbf{w}^I_t - \mathbf{w}^U_t\| + \frac{M_1r\eta}{n-r}(A(1 - \eta\mu)^{j_0+xT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}).$$

This can be used recursively until iteration  $j_m = j_0 + (x + m)T_0 - 1$ , i.e.:

$$\begin{aligned} &\leq (1 - \eta C)^{t-(j_0+(x+m-1)T_0)-1}\|\mathbf{w}^I_{j_0+(x+m-1)T_0+1} - \mathbf{w}^U_{j_0+(x+m-1)T_0+1}\| \\ &+ \frac{1 - (1 - \eta C)^{t-(j_0+(x+m-1)T_0)}}{\eta C} \frac{M_1r\eta}{n-r}(A(1 - \eta\mu)^{j_0+xT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}) \\ &\leq (1 - \eta C)^{t-(j_0+(x+m-1)T_0)-1}\|\mathbf{w}^I_{j_0+(x+m-1)T_0+1} - \mathbf{w}^U_{j_0+(x+m-1)T_0+1}\| \\ &+ \frac{M_1r}{C(n-r)}(A(1 - \eta\mu)^{j_0+xT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}) \end{aligned}$$

We can set  $t = j_0 + (y + m)T_0$  and for any  $y = 1, 2, \dots, x - 1$ , the formula above can be rewritten as:

$$\begin{aligned} &\|\mathbf{w}^I_{j_0+(y+m)T_0} - \mathbf{w}^U_{j_0+(y+m)T_0}\| \\ &\leq (1 - \eta C)^{T_0-1}\|\mathbf{w}^I_{j_0+(y+m-1)T_0+1} - \mathbf{w}^U_{j_0+(y+m-1)T_0+1}\| \\ &+ \frac{M_1r}{C(n-r)}(A(1 - \eta\mu)^{j_0+yT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}) \end{aligned}$$

Then at the iteration  $t = j_0 + (y + m - 1)T_0$ , the gradient is explicitly evaluated, which means that:

$$\|\mathbf{w}^I_{j_0+(y+m-1)T_0+1} - \mathbf{w}^U_{j_0+(y+m-1)T_0+1}\| \leq (1 - \eta\mu)\|\mathbf{w}^I_{j_0+(y+m-1)T_0} - \mathbf{w}^U_{j_0+(y+m-1)T_0}\|.$$

Since  $C = \mu - \frac{n}{n-r}\xi_{j_0, j_0+(m-1)T_0} - \frac{c_0M_1r}{2n}$ , then  $1 - \eta\mu < 1 - \eta C$  and thus

$$\|\mathbf{w}^I_{j_0+(y+m-1)T_0+1} - \mathbf{w}^U_{j_0+(y+m-1)T_0+1}\| \leq (1 - \eta C)\|\mathbf{w}^I_{j_0+(y+m-1)T_0} - \mathbf{w}^U_{j_0+(y+m-1)T_0}\|,$$

which can be plugged into the formula above:

$$\begin{aligned} &\|\mathbf{w}^I_{j_0+(y+m)T_0} - \mathbf{w}^U_{j_0+(y+m)T_0}\| \\ &\leq (1 - \eta C)^{T_0}\|\mathbf{w}^I_{j_0+(y+m-1)T_0} - \mathbf{w}^U_{j_0+(y+m-1)T_0}\| \\ &+ \frac{M_1r}{C(n-r)}(A(1 - \eta\mu)^{j_0+yT_0}d_{0, mT_0-1} + A\frac{1}{\frac{1}{2}-\frac{r}{n}}M_1\frac{r}{n}). \end{aligned}$$

This can be used recursively over  $y = x - 1, x - 2, \dots, 2, 1$ :

$$\begin{aligned}
& \|\mathbf{w}^I_{j_0+(y+m)T_0} - \mathbf{w}^U_{j_0+(y+m)T_0}\| \\
& \leq (1 - \eta C)^{yT_0} \|\mathbf{w}^I_{j_0+mT_0} - \mathbf{w}^U_{j_0+mT_0}\| \\
& + \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{M_1 r}{C(n-r)} (A(1 - \eta\mu)^{j_0+pT_0} d_{0,mT_0-1} + A \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n}) \\
& = (1 - \eta C)^{yT_0} \|\mathbf{w}^I_{j_0+mT_0} - \mathbf{w}^U_{j_0+mT_0}\| \\
& + \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{M_1 r}{C(n-r)} (A(1 - \eta\mu)^{j_0+pT_0} d_{0,mT_0-1}) \\
& + \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)},
\end{aligned} \tag{37}$$

in which

$$\begin{aligned}
& \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{M_1 r}{C(n-r)} (A(1 - \eta\mu)^{j_0+pT_0} d_{0,mT_0-1}) \\
& = \frac{AM_1 r \eta}{C(n-r)} (1 - \eta C)^{yT_0} (1 - \eta\mu)^{j_0} d_{0,mT_0-1} \sum_{p=1}^y (1 - \eta C)^{-pT_0} (1 - \eta\mu)^{pT_0}.
\end{aligned}$$

Recall that since  $1 - \eta C > 1 - \eta\mu$ , then the formula above can be bounded as:

$$\begin{aligned}
& \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{M_1 r}{C(n-r)} (A(1 - \eta\mu)^{j_0+pT_0} d_{0,mT_0-1}) \\
& \leq \frac{AM_1 r \eta}{C(n-r)} (1 - \eta C)^{yT_0} (1 - \eta\mu)^{j_0} d_{0,mT_0-1} \frac{1}{1 - (\frac{1-\eta\mu}{1-\eta C})^{T_0}}.
\end{aligned}$$

Also  $\sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)}$  can be simplified to:

$$\begin{aligned}
& \sum_{p=1}^y (1 - \eta C)^{(y-p)T_0} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)} = \sum_{p=0}^{y-1} (1 - \eta C)^{pT_0} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)} \\
& \leq \frac{1}{1 - (1 - \eta C)^{T_0}} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)}.
\end{aligned}$$

So equation (37) can be further bounded as:

$$\begin{aligned}
& \|\mathbf{w}^I_{j_0+(y+m)T_0} - \mathbf{w}^U_{j_0+(y+m)T_0}\| \\
& \leq (1 - \eta C)^{yT_0} \|\mathbf{w}^I_{j_0+mT_0} - \mathbf{w}^U_{j_0+mT_0}\| \\
& + \frac{AM_1 r}{C(n-r)} (1 - \eta C)^{yT_0} (1 - \eta\mu)^{j_0} d_{0,mT_0-1} \frac{1}{1 - (\frac{1-\eta\mu}{1-\eta C})^{T_0}} \\
& + \frac{1}{1 - (1 - \eta C)^{T_0}} \frac{AM_1^2 r^2}{C(n-r)(n/2-r)}.
\end{aligned} \tag{38}$$

When  $t \rightarrow \infty$  and thus  $y \rightarrow \infty$ ,  $(1 - \eta C)^{yT_0} \rightarrow 0$  and thus

$$\|\mathbf{w}^I_{j_0+(y+m)T_0} - \mathbf{w}^U_{j_0+(y+m)T_0}\| = o\left(\frac{r}{n}\right).$$

□

### A.3 Results for stochastic gradient descent

#### A.3.1 Quasi-Newton

We modify Equations (13) and (12) to SGD versions:

$$\mathbf{B}^S_{j_{k+1}} = \mathbf{B}^S_{j_k} - \frac{\mathbf{B}^S_{j_k} \Delta w^S_{j_k} \Delta w^{S^T}_{j_k} \mathbf{B}^S_{j_k}}{\Delta w^{S^T}_{j_k} \mathbf{B}^S_{j_k} \Delta w^S_{j_k}} + \frac{\Delta g^S_{j_k} \Delta g^{S^T}_{j_k}}{\Delta g^{S^T}_{j_k} \Delta w^S_{j_k}} \quad (39)$$

$$\mathbf{B}^{S^{-1}}_{j_{k+1}} = \left( \mathbf{I} - \frac{\Delta w^S_{j_k} \Delta g^{S^T}_{j_k}}{\Delta g^{S^T}_{j_k} \Delta w^S_{j_k}} \right) \mathbf{B}^{S^{-1}}_{j_k} \left( \mathbf{I} - \frac{\Delta g^S_{j_k} \Delta w^{S^T}_{j_k}}{\Delta g^{S^T}_{j_k} \Delta w^S_{j_k}} \right) + \frac{\Delta w^S_{j_k} \Delta w^{S^T}_{j_k}}{\Delta g^{S^T}_{j_k} \Delta w^S_{j_k}} \quad (40)$$

This iteration has the same initialization as  $\mathbf{B}_{j_k}$  and  $\mathbf{B}^{-1}_{j_k}$  but relies on the history information collected from the SGD-based training process  $[\Delta w^S_{j_0}, \Delta w^S_{j_1}, \dots, \Delta w^S_{j_{m-1}}]$  and  $[\Delta g^S_{j_0}, \Delta g^S_{j_1}, \dots, \Delta g^S_{j_{m-1}}]$  where  $\Delta w^S_{j_x} = \mathbf{w}^S_{j_x} - \mathbf{w}^{I,S}_{j_x}$  and  $\Delta g^S_{j_x} = G_{B,S}(\mathbf{w}^{I,S}_{j_x}) - G_{B,S}(\mathbf{w}^S_{j_x})$  ( $x = 0, 1, 2, \dots, m-1$ ). By the same argument as the proof of Lemma 6, the following inequality holds:

$$K_1 \|\mathbf{z}\|^2 \leq \mathbf{z}^T \mathbf{B}^S_{j_k} \mathbf{z} \leq K_2 \|\mathbf{z}\|^2 \quad (41)$$

where  $K_1 := \frac{1}{(1+\frac{L}{\mu})^{2m} \frac{L}{\mu} + \frac{1-(1+\frac{L}{\mu})^{2m}}{1-(1+\frac{L}{\mu})^2} \frac{1}{\mu}}$  and  $K_2 := (m+1)L$ , which are both positive values representing a

lower bound and an upper bound on the eigenvalues of  $\mathbf{B}^S_{j_k}$ .

#### A.3.2 Proof preliminaries

Similar to the argument for the GD-version of DeltaGrad, we can give an upper bound on  $\delta_{t,S}$ :

**Lemma 9** (Upper bound on  $\delta_{t,S}$ ). *Define  $\delta_{t,S} = G_{B,S}(\mathbf{w}^{U,S}_t) - G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t)$ . Then  $\|\delta_{t,S}\| \leq 2c_2 \frac{\Delta B_t}{B}$ . Moreover, with probability higher than  $1 - t \times 2 \exp(-2\sqrt{B})$ ,*

$$\|\delta_{t',S}\| \leq 2c_2 \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$$

uniformly over all iterations  $t' \leq t$ .

*Proof.* Recall that

$$G_{B,S}(\mathbf{w}^{U,S}_t) = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^{U,S}_t)$$

and

$$G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t) = \frac{1}{B - \Delta B_t} \sum_{i \in \mathcal{B}_t, i \notin \mathcal{R}} \nabla F_i(\mathbf{w}^{U,S}_t).$$

By subtracting  $G_{B,S}(\mathbf{w}^{U,S}_t)$  from  $G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t)$ , we have:

$$\begin{aligned} & \|G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t) - G_{B,S}(\mathbf{w}^{U,S}_t)\| \\ &= \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^{U,S}_t) - \frac{1}{B - \Delta B_t} \sum_{i \in \mathcal{B}_t, i \notin \mathcal{R}} \nabla F_i(\mathbf{w}^{U,S}_t) \right\| \\ &= \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t, i \in \mathcal{R}} \nabla F_i(\mathbf{w}^{U,S}_t) + \left( \frac{1}{B} - \frac{1}{B - \Delta B_t} \right) \sum_{i \in \mathcal{B}_t, i \notin \mathcal{R}} \nabla F_i(\mathbf{w}^{U,S}_t) \right\| \end{aligned}$$

Then by using the triangle inequality and the fact that  $\|\nabla F_i(\mathbf{w}^{U,S}_t)\| \leq c_2$  (Assumption 3), the formula above can be bounded by  $\frac{2\Delta B_t c_2}{B}$ .

Because of the randomness from SGD, the  $r$  removed samples can be viewed as uniformly distributed among all  $n$  training samples. Each sample is included in a mini-batch according to the outcome of a Bernoulli( $r/n$ ) random variable  $\mathbf{S}_i$ . Within a single mini-batch  $\mathcal{B}_{t'}$  at the iteration  $t'$ , we get  $\mathbb{E}(\sum_{i \in \mathcal{B}_{t'}} \mathbf{S}_i) =$

$\mathbb{E}(\Delta B_{t'}) = B \frac{r}{n}$  and  $\text{Var}(\sum_{i \in \mathcal{B}_{t'}} \mathbf{S}_i) = B \frac{r}{n} (1 - \frac{r}{n})$ . So in terms of the random variable  $\frac{\Delta B_{t'}}{B}$ , its expectation and variance will be  $\mathbb{E}(\frac{\Delta B_{t'}}{B}) = \frac{r}{n}$  and  $\text{Var}(\frac{\Delta B_{t'}}{B}) = \frac{r}{n} (1 - \frac{r}{n})$ .

Then based on Hoeffding's inequality, the following inequality holds:

$$\Pr\left(\left|\frac{\Delta B_{t'}}{B} - \frac{r}{n}\right| \leq \epsilon\right) \geq 1 - 2 \exp(-2\epsilon^2 B).$$

Then by setting  $\epsilon = \frac{1}{B^{1/4}}$  the formula above can be written as:

$$\Pr\left(\left|\frac{\Delta B_{t'}}{B} - \frac{r}{n}\right| \geq \frac{1}{B^{1/4}}\right) \leq 2 \exp(-2\sqrt{B})$$

Then by taking the union for all the iterations before  $t$ , we get: with probability higher than  $1 - t \times 2 \exp(-2\sqrt{B})$ ,

$$\left|\frac{\Delta B_{t'}}{B} - \frac{r}{n}\right| \leq \frac{1}{B^{1/4}}$$

and thus

$$\frac{\Delta B_{t'}}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}} \quad (42)$$

for all  $t' \leq t$ . □

In what follows, we use  $\Psi_1$  to represent  $\Psi_1 := 2 \exp(-2\sqrt{B})$ , which goes to 0 with large  $B$ .

Next we provide a bound for the sum of random sampled Hessian matrices within a minibatch in SGD.

**Theorem 6** (Hessian matrix bound in SGD). *With probability higher than*

$$1 - 2p \exp\left(-\frac{\log(2p)\sqrt{B}}{4 + \frac{2}{3}\left(\frac{\log^2(2p)}{B}\right)^{1/4}}\right),$$

for a given iteration  $t$ ,  $\left\|\left(\frac{1}{B} \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w}^{S_t})\right) - \mathbf{H}(\mathbf{w}^{S_t})\right\| \leq L \left(\frac{\log^2(2p)}{B}\right)^{1/4}$  where  $p$  represents the number of model parameters.

*Proof.* We consider using the matrix Bernstein inequality, Lemma 5. We define the random matrix  $\mathbf{S}_i = \frac{\mathbf{H}_i(\mathbf{w}) - \mathbf{H}(\mathbf{w})}{B}$  ( $i \in \mathcal{B}_t$ ). Due to the randomness from SGD, we know that  $\mathbb{E}(\mathbf{S}_i) = \mathbb{E}\left(\frac{\mathbf{H}_i(\mathbf{w}) - \mathbf{H}(\mathbf{w})}{B}\right) = \mathbf{0}$ . Using the sum  $\mathbf{Z}$  as required in Lemma 5,  $\mathbf{Z} = \left(\frac{1}{B} \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w})\right) - \mathbf{H}(\mathbf{w})$ . Also note that  $\mathbf{H}(\mathbf{w})$  and  $\mathbf{H}_i(\mathbf{w})$  are both  $p \times p$  matrices, so  $d_1 = d_2 = p$  in Lemma 5.

Furthermore, for each  $\mathbf{S}_i = \frac{\mathbf{H}_i(\mathbf{w}) - \mathbf{H}(\mathbf{w})}{B}$ , its norm is bounded by  $\frac{2L}{B}$  based on the smoothness condition, which means that  $J = \frac{2L}{B}$  in Lemma 5. Then we can explicitly calculate the upper bound on  $E(\mathbf{S}_i \mathbf{S}_i^*)$  and  $V(\mathbf{Z})$ :

$$\begin{aligned} \|\mathbb{E}(\mathbf{S}_i \mathbf{S}_i^*)\| &\leq \mathbb{E}(\|\mathbf{S}_i \mathbf{S}_i^*\|) \leq \mathbb{E}(\|\mathbf{S}_i\| \|\mathbf{S}_i^*\|) \leq J^2 = \frac{4L^2}{B^2}, \\ V(\mathbf{Z}) &\leq \sum_{i \in \mathcal{B}_t} \frac{4L^2}{B^2} = \frac{4L^2}{B}. \end{aligned}$$

Thus by plugging the above expression into equation (9) and (10), we get:

$$\begin{aligned} P(\|\mathbf{Z}\| \geq x) &= \Pr\left(\left\|\left(\frac{1}{B} \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w})\right) - \mathbf{H}(\mathbf{w})\right\| \geq x\right) \\ &\leq (d_1 + d_2) \exp\left(\frac{-x^2}{\frac{4L^2}{B} + \frac{2Lx}{3B}}\right) = 2p \exp\left(\frac{-x^2}{\frac{4L^2}{B} + \frac{2Lx}{3B}}\right), \forall x \geq 0 \end{aligned} \quad (43)$$

$$\begin{aligned}
\mathbb{E}(\|\mathbf{Z}\|) &= \mathbb{E} \left( \left\| \left( \frac{1}{B} \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w}) \right) - \mathbf{H}(\mathbf{w}) \right\| \right) \\
&\leq \sqrt{\frac{8L^2}{B} \log(d_1 + d_2)} + \frac{2L}{3B} \log(d_1 + d_2) = \sqrt{\frac{8L^2}{B} \log(2p)} + \frac{2L}{3B} \log(2p).
\end{aligned} \tag{44}$$

Then by setting  $x = L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$ , Equation (43) becomes:

$$\begin{aligned}
&\Pr(\|\mathbf{Z}\| \geq L \left( \frac{\log^2(2p)}{B} \right)^{1/4}) \\
&= \Pr \left( \left\| \left( \frac{1}{B} \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w}) \right) - \mathbf{H}(\mathbf{w}) \right\| \geq L \left( \frac{\log^2(2p)}{B} \right)^{1/4} \right) \\
&\leq (2p) \exp \left( \frac{-\frac{L^2 \log(2p)}{\sqrt{B}}}{\frac{4L^2}{B} + \frac{2L^2}{3B} \left( \frac{\log^2(2p)}{B} \right)^{1/4}} \right) = (2p) \exp \left( -\frac{\log(2p)\sqrt{B}}{4 + \frac{2}{3} \left( \frac{\log^2(2p)}{B} \right)^{1/4}} \right).
\end{aligned} \tag{45}$$

For large mini-batch size  $B$ , both  $L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$  and  $(2p) \exp \left( -\frac{\log(2p)\sqrt{B}}{4 + \frac{2}{3} \left( \frac{\log^2(2p)}{B} \right)^{1/4}} \right)$  are approaching 0.  $\square$

In what follows, we use  $\Psi_2$  to denote the probability  $\Psi_2 := (2p) \exp \left( -\frac{\log(2p)\sqrt{B}}{4 + \frac{2}{3} \left( \frac{\log^2(2p)}{B} \right)^{1/4}} \right)$ .

Based on this result, we can derive an SGD version of Theorem 1 as below, which also relies on a preliminary estimate on the bound on  $\|\mathbf{w}^{I,S_t} - \mathbf{w}^{S_t}\|$ :

**Theorem 7** (Error in mean Hessian, and in secant equation with incorrect quasi-Hessian for SGD). *Suppose that  $\|\mathbf{w}^{S_{t'}} - \mathbf{w}^{I,S_{t'}}\| \leq M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$  and*

$$\left\| \left( \frac{1}{B} \sum_{i \in \mathcal{B}_{t'}} \mathbf{H}_i(\mathbf{w}^{S_{t'}}) \right) - \mathbf{H}(\mathbf{w}^{S_{t'}}) \right\| \leq L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$$

hold for any  $t' \leq t$  where  $M_1 = \frac{2c_2}{\mu}$ ,  $\mu$  is from Assumption 2 and  $c_2$  is from Assumption 3. Let  $e = \frac{L(L+1)+K_2L}{\mu K_1}$  for the upper and lower bounds  $K_1, K_2$  on the eigenvalues of the quasi-Hessian from equation (41) and for the Lipschitz constant  $c_0$  of the Hessian. For any  $t_1, t_2$  such that  $1 \leq t_1 < t_2 \leq t$ , we have:

$$\|\mathbf{H}^{S_{t_1}} - \mathbf{H}^{S_{t_2}}\| \leq 2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + c_0 d_{t_1, t_2} + 3c_0 M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right).$$

For any  $j_1, j_2, \dots, j_m$  such that  $j_m \leq t' \leq j_m + T_0 - 1$  and  $t' \leq t$ , we have:

$$\begin{aligned}
&\|\Delta g_{j_k}^S - \mathbf{B}_{j_q}^S \Delta w_{j_k}^S\| \\
&\leq [(1+e)^{j_q - j_k - 1} - 1] \cdot [2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + c_0 d_{j_k, j_q} + \frac{3c_0 M_1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)] \cdot s_{j_m, j_1}.
\end{aligned}$$

Here  $s_{j_m, j_1} = \max(\|\Delta w^S\|)_{a=j_1, j_2, \dots, j_m}$ ,  $d_{j_k, j_q} = \max(\|\mathbf{w}^S_a - \mathbf{w}^S_b\|)_{j_k \leq a \leq b \leq j_q}$ ,  $\mathbf{H}^S_t$  is the average of the Hessian matrix evaluated between  $\mathbf{w}^{S_t}$  and  $\mathbf{w}^{I,S_t}$  for the samples in mini-batch  $\mathcal{B}_t$ :

$$\mathbf{H}^S_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \int_0^1 \mathbf{H}_i(\mathbf{w}^S_t + x(\mathbf{w}^{U,S_t} - \mathbf{w}^S_t)) dx.$$

*Proof.* First of all, let us bound  $\|\mathbf{H}^S_{t_1} - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx\|$  by adding and subtracting  $\frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1})$  and  $\mathbf{H}(\mathbf{w}^S_{t_1})$  inside the norm:

$$\begin{aligned}
& \|\mathbf{H}^S_{t_1} - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx\| \\
&= \left\| \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx \right\| \\
&= \left\| \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} (\mathbf{H}_i(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1})) - \mathbf{H}_i(\mathbf{w}^S_{t_1}))dx + \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1}) \right. \\
&\quad \left. - \int_0^1 (\mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1})) - \mathbf{H}(\mathbf{w}^S_{t_1}))dx - \mathbf{H}(\mathbf{w}^S_{t_1}) \right\|.
\end{aligned}$$

Then by using the triangle inequality and Assumption 4, the formula above can be bounded as:

$$\begin{aligned}
& \leq \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \|\mathbf{H}_i(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1})) - \mathbf{H}_i(\mathbf{w}^S_{t_1})\| dx \\
& \quad + \int_0^1 \|\mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1})) - \mathbf{H}(\mathbf{w}^S_{t_1})\| dx \\
& \quad + \left\| \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1}) - \mathbf{H}(\mathbf{w}^S_{t_1}) \right\| \\
& \leq \frac{1}{B} \left( \sum_{i \in \mathcal{B}_{t_1}} \int_0^1 c_0 x \|\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}\| dx \right) + \int_0^1 c_0 x \|\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}\| dx \\
& \quad + \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1}) - \mathbf{H}(\mathbf{w}^S_{t_1}) \right\| \\
& \leq c_0 \|\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}\| + \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1}) - \mathbf{H}(\mathbf{w}^S_{t_1}) \right\|.
\end{aligned} \tag{46}$$

Then based on the above results, we can compute the bound on  $\|\mathbf{H}^S_{t_1} - \mathbf{H}^S_{t_2}\|$ , for which we use the triangle inequality first:

$$\begin{aligned}
& \|\mathbf{H}^S_{t_1} - \mathbf{H}^S_{t_2}\| \\
&= \left\| \mathbf{H}^S_{t_1} - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx \right\| \\
& \quad + \left\| \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_2} + x(\mathbf{w}^{I,S}_{t_2} - \mathbf{w}^S_{t_2}))dx \right\| \\
& \quad + \left\| \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_2} + x(\mathbf{w}^{I,S}_{t_2} - \mathbf{w}^S_{t_2}))dx - \mathbf{H}^S_{t_2} \right\|.
\end{aligned} \tag{47}$$

Then by using the result from Formula (46), this term can be further bounded as:

$$\begin{aligned}
& \leq c_0 \|\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}\| + \left\| \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_1}} \mathbf{H}_i(\mathbf{w}^S_{t_1}) - \mathbf{H}(\mathbf{w}^S_{t_1}) \right\| \\
& \quad + c_0 \|\mathbf{w}^{I,S}_{t_2} - \mathbf{w}^S_{t_2}\| + \left\| \int_0^1 \frac{1}{B} \sum_{i \in \mathcal{B}_{t_2}} \mathbf{H}_i(\mathbf{w}^S_{t_2}) - \mathbf{H}(\mathbf{w}^S_{t_2}) \right\| \\
& \quad + \left\| \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_1} + x(\mathbf{w}^{I,S}_{t_1} - \mathbf{w}^S_{t_1}))dx - \int_0^1 \mathbf{H}(\mathbf{w}^S_{t_2} + x(\mathbf{w}^{I,S}_{t_2} - \mathbf{w}^S_{t_2}))dx \right\|.
\end{aligned}$$

Since

$$\left\| \left( \frac{1}{B} \sum_{i \in \mathcal{B}_{t'}} \mathbf{H}_i(\mathbf{w}^{S_{t'}}) \right) - \mathbf{H}(\mathbf{w}^{S_{t'}}) \right\| \leq L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$$

for any  $t' \leq t$ , then the formula above can be bounded as:

$$\begin{aligned} &\leq 2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + c_0 \|\mathbf{w}^{S_{t_1}} - \mathbf{w}^{S_{t_2}}\| + \frac{c_0}{2} \|\mathbf{w}^{S_{t_1}} - \mathbf{w}^{I,S_{t_1}}\| \\ &+ \frac{c_0}{2} \|\mathbf{w}^{I,S_{t_2}} - \mathbf{w}^{S_{t_2}}\| + c_0 \|\mathbf{w}^{I,S_{t_1}} - \mathbf{w}^{S_{t_1}}\| + c_0 \|\mathbf{w}^{I,S_{t_2}} - \mathbf{w}^{S_{t_2}}\| \\ &= 2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + 3c_0 M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + c_0 d_{t_1, t_2}. \end{aligned}$$

This finishes the proof of the first inequality. Then by defining

$$f = \left( 2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + 3c_0 M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + c_0 d_{j_k, j_q} \right) s_{j_m, j_1}$$

and using the same argument as Equation (21)-(23) (except that  $\Delta w$  and  $\Delta g$  are replaced with  $\Delta w^S$  and  $\Delta g^S$ ), the following inequality thus holds:

$$\begin{aligned} b_{j_q} &= \left\| \Delta g^S_{j_k} - \left( \mathbf{B}^S_{j_q} - \frac{\mathbf{B}^S_{j_q} \Delta w^S_{j_q} \Delta w^{S^T}_{j_q} \mathbf{B}^S_{j_q}}{\Delta w^{S^T}_{j_q} \mathbf{B}^S_{j_q} \Delta w^S_{j_q}} + \frac{\Delta g^S_{j_q} \Delta g^{S^T}_{j_q}}{\Delta g^{S^T}_{j_q} \Delta w^S_{j_q}} \right) \Delta w^S_{j_k} \right\| \\ &\leq [(1+e)^{j_q - j_k} - 1] f \end{aligned} \quad (48)$$

and thus

$$\|\Delta g^S_{j_k} - \mathbf{B}^S_{j_q} \Delta w^S_{j_k}\| \leq [(1+e)^{j_q - j_k - 1} - 1] f,$$

which finishes the proof.

For simplicity, we denote  $M_1^S := M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}}$ . So the preliminary estimate of the bound on  $\|\mathbf{w}^{S_{t'}} - \mathbf{w}^{I,S_{t'}}\|$  becomes:  $\|\mathbf{w}^{S_{t'}} - \mathbf{w}^{I,S_{t'}}\| \leq M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$   $\square$

Similarly, we get a SGD-version of Corollary 1:

**Corollary 2** (Approximation accuracy of Quasi-Hessian to mean Hessian). *Suppose that  $\|\mathbf{w}^{S_{t'}} - \mathbf{w}^{I,S_{t'}}\| \leq M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$  and*

$$\left\| \left( \frac{1}{B} \sum_{i \in \mathcal{B}_{t'}} \mathbf{H}_i(\mathbf{w}^{S_{t'}}) \right) - \mathbf{H}(\mathbf{w}^{S_{t'}}) \right\| \leq L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$$

*hold for any  $t' \leq t$ .  $M_1$  and  $M_1^S$  are provided in Theorem 7, i.e.  $M_1 = \frac{2c_2}{\mu}$  and  $M_1^S = M_1 \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}}$ . Then for any  $t'$  and  $j_m$  such that  $j_m \leq t' \leq j_m + T_0 - 1$  and  $t' \leq t$ , the following inequality holds:*

$$\begin{aligned} \|\mathbf{H}^{S_{t'}} - \mathbf{B}^S_{j_m}\| &\leq \xi_{j_1, j_m}^S \\ &:= A(d_{j_1, j_m + T_0 - 1} + 3M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + \frac{2}{c_0} L \left( \frac{\log^2(2p)}{B} \right)^{1/4}) \end{aligned}$$

*where recall again that  $c_0$  is the Lipschitz constant of the Hessian,  $d_{j_1, j_m + T_0 - 1}$  is the maximal gap between the iterates of the SGD algorithm on the full data from  $j_1$  to  $j_m + T_0 - 1$  and  $A = \frac{c_0 \sqrt{m} [(1+e)^m - 1]}{c_1} + c_0$  in which  $e$  is a problem dependent constant defined in Theorem 7,  $c_1$  is the "strong independence" constant from Assumption (5).*



This proof is similar to the proof of Corollary 1. First of all,  $\mathbf{H}$ ,  $\mathbf{B}$ ,  $\xi_{j_1, j_m}$  in Corollary 1 are replaced with  $\mathbf{H}^S$ ,  $\mathbf{B}^S$ ,  $\xi_{j_1, j_m}^S$ . Second, Theorem 7 holds and thus the following inequality holds:

$$\|\mathbf{H}^S_{t'} - \mathbf{H}^S_{j_m}\| \leq 2L \left( \frac{\log^2(2p)}{B} \right)^{1/4} + c_0 d_{t', j_m} + 3c_0 M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$$

by using strong independence from Assumption 5,  $\|\mathbf{H}^S_{j_m} - \mathbf{B}^S_{j_m}\|$  can be bounded as:

$$\|\mathbf{H}^S_{j_m} - \mathbf{B}^S_{j_m}\| \leq \sqrt{m} [(1+e)^m - 1] \frac{c_0}{c_1} \cdot \left( d_{j_1, j_m + T_0 - 1} + 3M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + \frac{2}{c_0} L \left( \frac{\log^2(2p)}{B} \right)^{1/4} \right) \quad (49)$$

Then by combining the two formulas above, we know that Corollary 2 holds. Note that the definition of  $\xi_{j_1, j_m}^S$  can be rewritten as below:

$$\begin{aligned} \xi_{j_1, j_m}^S &= A(d_{j_1, j_m + T_0 - 1} + 3M_1^S \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + \frac{2}{c_0} L \left( \frac{\log^2(2p)}{B} \right)^{1/4}) \\ &=: A d_{j_1, j_m + T_0 - 1} + A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} \end{aligned} \quad (50)$$

in which  $A_1 := 3AM_1^S$  and  $A_2 := 3AM_1^S + \frac{2AL(\log(2p))^{1/2}}{c_0}$ .

We can do a similar analysis to Lemma 8 by simply replacing  $\mathbf{w}_t$  and  $F(\ast)$  with  $\mathbf{w}_t^S$  and  $G_{B,S}$ :

**Lemma 10.** *Let us use the definition of  $d_{k,q}$  from Theorem 7:*

$$d_{k,q} = \max(\|\mathbf{w}_a^S - \mathbf{w}_b^S\|)_{k \leq a \leq b \leq q}$$

where  $k < q \leq t$ , then  $d_{k,q} \leq (1 - \eta\mu)^k d_{0,q-j} + 2c_2 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}$  holds with probability higher than  $1 - t(p+1) \exp\left(-\frac{\log(p+1)\sqrt{B}}{4 + \frac{2}{3} \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}}\right)$ .

*Proof.* According to Lemma 5, we can define a random matrix  $\mathbf{S}_i = \frac{1}{B}(\nabla F_i(\mathbf{w}^S_a) - \nabla F(\mathbf{w}^S_a))$  where recall that  $\nabla F(\mathbf{w}^S_a) = \frac{1}{n} \sum_{i=1}^n \nabla F_{i,a}(\mathbf{w}^S_a)$  ( $i \in \mathcal{B}_t$ ). Due to the randomness from SGD, we know that  $\mathbb{E}(\mathbf{S}_i) = \mathbf{0}$ . Based on the definition of  $\mathbf{Z}$  in Lemma 5,  $\mathbf{Z} = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_a) - \nabla F(\mathbf{w}^S_a)$ . Also note that  $\nabla F_i(\mathbf{w}^S_a)$  and  $\nabla F(\mathbf{w}^S_a)$  are both  $p \times 1$  matrices, so  $d_1 = p$  and  $d_2 = 1$  in Lemma 5.

Moreover according to Assumption 3,  $\|\nabla F_i(\mathbf{w}^S_a)\| \leq c_2$ . Then we know that  $V(\mathbf{Z}) \leq \frac{4c_2^2}{B}$  and  $\|\mathbf{S}_i\| \leq \frac{2c_2}{B}$ . So according to Lemma 5, the following inequality holds:

$$\begin{aligned} P(\|\mathbf{Z}\| \geq x) &= \Pr\left(\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_a) - \nabla F(\mathbf{w}^S_a) \right\| \geq x\right) \\ &\leq (d_1 + d_2) \exp\left(\frac{-x^2}{\frac{4c_2^2}{B} + \frac{2c_2x}{3B}}\right) = (p+1) \exp\left(\frac{-x^2}{\frac{4c_2^2}{B} + \frac{2c_2x}{3B}}\right), \forall x \geq 0 \end{aligned} \quad (51)$$

By setting  $x = c_2 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}$ , the formula above is evaluated as:

$$\begin{aligned} \Pr\left(\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}^S_a) - \nabla F(\mathbf{w}^S_a) \right\| \geq c_2 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}\right) \\ \leq (p+1) \exp\left(\frac{-\log(p+1)\sqrt{B}}{4 + \frac{2}{3} \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}}\right) \end{aligned}$$

So by taking the union for the first  $t$  iterations, then with probability higher than  $1 - t(p+1) \exp\left(-\frac{\log(p+1)\sqrt{B}}{4 + \frac{2}{3}\left(\frac{(\log(p+1))^2}{B}\right)^{1/4}}\right)$ , the following inequality holds for all  $t' \leq t$ :

$$\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_{t'}} \nabla F_i(\mathbf{w}^S_{a_i}) - \nabla F(\mathbf{w}^S_a) \right\| \leq c_2 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4} \quad (52)$$

Then by using the similar arguments to Lemma 8, we get:

$$\|\mathbf{w}^S_a - \mathbf{w}^S_b\| \leq (1 - \eta\mu)^z \|\mathbf{w}^S_{a-z} - \mathbf{w}^S_{b-z}\| + \frac{2c_2}{\mu} \left( \frac{(\log(p+1))^2}{B} \right)^{1/4} = (1 - \eta\mu)^z \|\mathbf{w}^S_{a-z} - \mathbf{w}^S_{b-z}\| + M_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}$$

and thus  $d_{k,q} \leq (1 - \eta\mu)^k d_{0,q-k} + M_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}$  holds with probability higher than  $1 - t(p+1) \exp\left(-\frac{\log(p+1)\sqrt{B}}{4 + \frac{2}{3}\left(\frac{(\log(p+1))^2}{B}\right)^{1/4}}\right)$ .

In what follows, we use  $\Psi_3$  to denote  $(p+1) \exp\left(-\frac{\log(p+1)\sqrt{B}}{4 + \frac{2}{3}\left(\frac{(\log(p+1))^2}{B}\right)^{1/4}}\right)$ .

□

Then by using the definition of  $\xi_{j_1, j_m}^S$ , the following inequality holds with probability higher than  $1 - t\Psi_3$  for any  $x$  such that for  $j_0 + (x + m - 1)T_0 \leq t$ , the following inequality holds:

$$\begin{aligned} \xi_{j_1, j_m}^S &= \xi_{j_0 + xT_0, j_0 + (x+m-1)T_0}^S \leq (1 - \eta\mu)^{xT_0} Ad_{j_0, j_0 + mT_0 - 1} \\ &+ A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4} \end{aligned} \quad (53)$$

### A.3.3 Main recursions

We bound the difference between  $\mathbf{w}^{I,S}_t$  and  $\mathbf{w}^{U,S}_t$ . First we bound  $\|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t\|$ :

**Theorem 8** (Bound between iterates on full and the leave- $r$ -out dataset). *When*

$$\frac{\Delta B_{t'}}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}}$$

holds for all  $t' < t$ ,  $\|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t\| \leq \frac{2c_2}{\mu} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$ . Since with probability higher than  $1 - t \times \Psi_1$ ,

$$\frac{\Delta B_{t'}}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}}$$

holds for all  $t' < t$ . Then with the same probability,  $\|\mathbf{w}^S_{t'+1} - \mathbf{w}^{U,S}_{t'+1}\| \leq M_1 \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right)$  for all iterations  $t' < t$ , where recall that  $M_1 = \frac{2c_2}{\mu}$ .

Similarly, we can bound the difference between  $\mathbf{w}^I_t$  and  $\mathbf{w}_t$ .

**Theorem 9** (Bound between iterates on full data and incrementally updated ones). *Suppose that for at some iteration  $t$  and any given  $t' \leq t$  such that  $j'_m \leq t' \leq j'_m + T_0 - 1$ , we have the following bounds:*

1.  $\|\mathbf{H}^S_{t'} - \mathbf{B}^S_{j'_m}\| \leq \xi_{j'_1, j'_m}^S = Ad_{j'_1, j'_m + T_0 - 1} + A \frac{3}{\frac{1}{2} - \frac{r}{n}} M_1 \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) + A \frac{2}{c_0} L \left( \frac{\log^2(2p)}{B} \right)^{1/4}$ ;
2.  $\frac{\Delta B_{t'}}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}}$ ;
3. Formula (53) holds for any  $x$  such that  $j_0 + (x + m - 1)T_0 \leq t$ ;
4.  $\xi_{j_0, j_0 + (m-1)T_0}^S + A \times M_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4} \leq \frac{\mu}{2}$ ,

then

$$\|\mathbf{w}^{I,S}_{t'+1} - \mathbf{w}^S_{t'+1}\| \leq \frac{2c_2}{(\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}})\mu} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right) = M_1^S \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$$

for any  $t' \leq t$ . Recall that  $c_0$  is the Lipschitz constant of the Hessian,  $M_1$  and  $A$  are defined in Theorem 8 and Corollary 2 respectively, and do not depend on  $t$ .

in particular for all  $t$ , the following inequality holds:

$$\|\mathbf{w}^{I,S}_{t+1} - \mathbf{w}^S_{t+1}\| \leq \frac{2c_2}{(\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}})\mu} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right) = M_1^S \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right).$$

Similarly, we will show that both inequalities  $\|\mathbf{H}^S_t - \mathbf{B}^S_{j_m}\| \leq \xi_{j_1, j_m}^S$  and  $\|\mathbf{w}^{I,S}_{t+1} - \mathbf{w}^S_{t+1}\| \leq M_1^S \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$  hold for all iterations  $t$ .

**Theorem 10** (Bound between iterates on full data and incrementally updated ones (all iterations)). *Suppose that there are  $T$  iterations in total for each training phase, then with probability higher than  $1 - T \times (\Psi_1 + \Psi_2 + \Psi_3)$ , for any  $t$  where  $j_m < t < j_m + T_0 - 1$ ,  $\|\mathbf{w}^{I,S}_t - \mathbf{w}^S_t\| \leq \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$  and  $\|\mathbf{H}^S_t - \mathbf{B}^S_{j_m}\| \leq \xi_{j_1, j_m}^S$ , where  $\xi_{j_1, j_m}^S$  is defined in Corollary 2,  $\Psi_1$  is defined in Lemma 9,  $\Psi_2$  is defined in Theorem 6 and  $\Psi_3$  is defined in Lemma 10.*

Then we have the following bound for  $\|\mathbf{w}^U_t - \mathbf{w}^I_t\|$ .

**Theorem 11** (Main result: Bound between true and incrementally updated iterates for SGD). *Suppose that there are  $T$  iterations in total for each training phase, then with probability higher than  $1 - T \times (\Psi_1 + \Psi_2 + \Psi_3)$ , the result  $\mathbf{w}^{I,S}_t$  of Algorithm 1 approximates the correct iteration values  $\mathbf{w}^{U,S}_t$  at the rate*

$$\|\mathbf{w}^{U,S}_t - \mathbf{w}^{I,S}_t\| \leq o\left(\frac{r}{n} + \frac{1}{B^{1/4}}\right).$$

So  $\|\mathbf{w}^{U,S}_t - \mathbf{w}^{I,S}_t\|$  is of a lower order than  $\left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$ .

### A.3.4 Proof of Theorem 8

*Proof.* By subtracting  $\mathbf{w}^S_t - \mathbf{w}^{U,S}_t$ , taking the matrix norm and using the update rule in equation (5) and (6), we get:

$$\begin{aligned} & \|\mathbf{w}^S_{t+1} - \mathbf{w}^{U,S}_{t+1}\| \\ &= \|\mathbf{w}^S_t - \eta G_{B,S}(\mathbf{w}^S_t) - (\mathbf{w}^{U,S}_t - \eta G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t))\| \\ &= \|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t - \eta(G_{B,S}(\mathbf{w}^S_t) - G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t))\| \\ &= \|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t - \eta(G_{B,S}(\mathbf{w}^S_t) - G_{B,S}(\mathbf{w}^{U,S}_t)) \\ &\quad + G_{B,S}(\mathbf{w}^{U,S}_t) - G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t))\| \\ &= \|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t - \eta(G_{B,S}(\mathbf{w}^S_t) - G_{B,S}(\mathbf{w}^{U,S}_t)) + \\ &\quad \eta(G_{B,S}(\mathbf{w}^{U,S}_t) - G_{B-\Delta B,S}^U(\mathbf{w}^{U,S}_t))\| \end{aligned} \tag{54}$$

By Cauchy mean-value theorem and the triangle inequality, the above formula becomes:

$$\begin{aligned} & \leq \|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t - \eta\left(\frac{1}{B} \int_0^1 \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w}^S_t + x(\mathbf{w}^{U,S}_t - \mathbf{w}^S_t)) dx\right)(\mathbf{w}^S_t - \mathbf{w}^{U,S}_t)\| + \eta \|\delta_{t,S}\| \\ &= \left\| \left( \mathbf{I} - \eta \left( \frac{1}{B} \int_0^1 \sum_{i \in \mathcal{B}_t} \mathbf{H}_i(\mathbf{w}^S_t + x(\mathbf{w}^{U,S}_t - \mathbf{w}^S_t)) dx \right) \right) (\mathbf{w}^S_t - \mathbf{w}^{U,S}_t) \right\| + \eta \|\delta_{t,S}\| \end{aligned}$$

Then by using the Lemma 3 and Lemma 9, the formula above can be bounded as:

$$\leq (1 - \eta\mu) \|\mathbf{w}^S_t - \mathbf{w}^{U,S}_t\| + \eta 2c_2 \frac{\Delta B_t}{B}$$

Then by using Lemma 9 and using the formula above recursively, we get that with probability higher than  $1 - t \cdot \Psi_1$ ,  $\|\mathbf{w}^{S_{t'}} - \mathbf{w}^{U, S_{t'}}\| \leq \frac{2c_2}{\mu} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$  holds for all iterations  $t' \leq t$ , which finishes the proof.  $\square$

### A.3.5 Proof of Theorem 9

*Proof.* For any  $t' \leq t$ , by subtracting  $\mathbf{w}^{S_{t'}}$  by  $\mathbf{w}^{I, S_{t'}}$  and taking the same argument as equation (29)-(31) (except that  $\mathbf{w}_{t'}$ ,  $\mathbf{w}^{I_{t'}}$ ,  $\mathbf{H}$ ,  $\mathbf{B}$ ,  $n$ ,  $r$  are replaced with  $\mathbf{w}^{S_{t'}}$ ,  $\mathbf{w}^{I, S_{t'}}$ ,  $\mathbf{H}^S$ ,  $\mathbf{B}^S$ ,  $B$ ,  $\Delta B_{t'}$ ), the following equality holds due to the bound on  $\|\mathbf{H}^{S_{t'}} - \mathbf{B}^{S_{j_m}}\|$ :

$$\begin{aligned} & \|\mathbf{w}^{I, S_{t'+1}} - \mathbf{w}^{S_{t'+1}}\| \\ & \leq (1 - \eta\mu + \eta \frac{B}{B - \Delta B_t} \xi_{j_1, j_m}^S) \|\mathbf{w}^{I, S_t} - \mathbf{w}^{S_t}\| + \frac{2\Delta B_t \eta c_2}{B - \Delta B_t}. \end{aligned} \quad (55)$$

Since  $\frac{\Delta B_t}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}}$  for all iterations between 0 and  $t$ , the following two inequalities hold:

$$\frac{2\Delta B_t \eta c_2}{B - \Delta B_t} = \frac{2\eta c_2}{\frac{B}{\Delta B_t} - 1} \leq \frac{2\eta c_2}{\frac{1}{\frac{r}{n} + \frac{1}{B^{1/4}}} - 1} = \frac{2\eta c_2}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right), \quad (56)$$

$$\frac{B}{B - \Delta B_t} = \frac{1}{1 - \frac{\Delta B_t}{B}} \leq \frac{1}{1 - \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)}. \quad (57)$$

Moreover, since Formula (53) holds and  $\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4} \leq \frac{\mu}{2}$ , then:

$$\begin{aligned} \xi_{j_1, j_m}^S &= \xi_{j_0+xT_0, j_0+(x+m-1)T_0}^S \leq (1 - \eta\mu)^{xT_0} Ad_{j_0, j_0+mT_0-1} \\ &+ A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4} \\ &\leq Ad_{j_0, j_0+mT_0-1} + A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4} \\ &= \xi_{j_0, j_0+mT_0-1} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4} \leq \frac{\mu}{2}. \end{aligned}$$

Then the Formula (55) can be bounded as:

$$\begin{aligned} & \|\mathbf{w}^{I, S_{t'+1}} - \mathbf{w}^{S_{t'+1}}\| \\ & \leq (1 - \eta\mu + \eta \frac{B}{B - \Delta B_t} (\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4})) \|\mathbf{w}^{I, S_{t'}} - \mathbf{w}^{S_{t'}}\| + \frac{2\Delta B_t \eta c_2}{B - \Delta B_t} \\ & \leq (1 - \eta\mu + \eta \frac{\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}}{1 - \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)}) \|\mathbf{w}^{I, S_{t'}} - \mathbf{w}^{S_{t'}}\| + \frac{2\eta c_2}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right), \end{aligned}$$

which uses equation (56) and (57). Then applying the formula recursively from iteration  $t$  to 0, we can get:

$$\begin{aligned} & \|\mathbf{w}^{I, S_{t'+1}} - \mathbf{w}^{S_{t'+1}}\| \\ & \leq \frac{1}{\eta\left(\mu - \frac{\xi_{j_0, j_0+(m-1)T_0}^S + 2c_2 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}}{1 - \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)}\right)} \frac{2\eta c_2}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right). \end{aligned}$$

Then since  $\xi_{j_0, j_0+(m-1)T_0}^S \leq \frac{\mu}{2}$ , the formula above can be further bounded as:

$$\begin{aligned} &= \frac{2c_2}{\left(1 - \frac{r}{n} - \frac{1}{B^{1/4}}\right)\mu - \xi_{j_0, j_0+(m-1)T_0}^S - A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right) \\ &\leq \frac{2c_2}{\left(\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}\right)\mu} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right). \end{aligned}$$

□

### A.3.6 Proof of Theorem 10

The proof is the same as the proof of Theorem 4 except that  $\mathbf{w}$ ,  $\mathbf{w}^I$ ,  $n$ ,  $r$ ,  $\xi_{j_1, j_m}$ ,  $\mathbf{H}$ ,  $\mathbf{B}$  need to be replaced by  $\mathbf{w}^S$ ,  $\mathbf{w}^{I,S}$ ,  $B$ ,  $r$ ,  $\xi_{j_1, j_m}^S$ ,  $\mathbf{H}^S$ ,  $\mathbf{B}^S$  and the main theorems that the proof depends on will be replaced by Theorem 9 and Corollary 2. But we need some careful explanations for the probability, which is shown as:

*Proof.* We define the following event at a given iteration  $k$ :

$$\begin{aligned} \Omega_0(k) &= \{\|\mathbf{w}_k^S - \mathbf{w}^{U,S}_k\| \leq M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)\}, \\ \Omega_1(k) &= \{\|\mathbf{w}_k^S - \mathbf{w}^{I,S}_k\| \leq M_1^S \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)\}, \\ \Omega_2(k) &= \{\|\mathbf{H}_{k-1}^S - \mathbf{B}_{j_m}^S\| \leq \xi_{j_1, j_m}^S\} \quad (j_m \leq k-1 \leq j_m + T_0 - 1), \\ \Omega_3(k) &= \left\{ \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_{k-1}} \mathbf{H}_i(\mathbf{w}_{k-1}^S) \right\| - \mathbf{H}(\mathbf{w}_{k-1}^S) \right\| \leq L \left(\frac{\log^2(2p)}{B}\right)^{1/4} \right\}, \\ \Omega_4(k) &= \{\xi_{j_0+xT_0, j_0+(x+m-1)T_0}^S \leq (1 - \eta\mu)^{j_0+xT_0} Ad_{0,mT_0-1} \\ &\quad + A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}\} \text{ where } j_0 + (x+m-1)T_0 \leq k-1 \leq j_0 + (x+m)T_0 - 1, \\ \Omega_5(k) &= \left\{ \frac{\Delta B_{k-1}}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}} \right\}. \end{aligned}$$

For all  $t$ , according to Corollary 2, the following equation holds:

$$\Pr\left(\bigcap_{k=1}^t \Omega_2(k) \mid \bigcap_{k=1}^{t-1} \Omega_1(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1.$$

in which the co-occurrence of multiple events is denoted by  $\bigcap$  or “,”. So this formula means that the probability that  $\Omega_2(k)$  is true for all  $k \leq t$  given that the events  $\Omega_1(k)$  and  $\Omega_3(k)$  are true at the same time for all  $k \leq t$  is 1.

Similarly, according to Theorem 9,  $\Pr\left(\bigcap_{k=1}^t \Omega_1(k) \mid \bigcap_{k=1}^t \Omega_2(k), \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k)\right) = 1$ . Then we know that:

$$\begin{aligned} &\Pr\left(\bigcap_{k=1}^t \Omega_1(k) \mid \bigcap_{k=1}^t \Omega_2(k), \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k)\right) \cdot \Pr\left(\bigcap_{k=1}^t \Omega_2(k) \mid \bigcap_{k=1}^{t-1} \Omega_1(k), \bigcap_{k=1}^t \Omega_3(k)\right) \\ &= \Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k) \mid \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^{t-1} \Omega_1(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1, \end{aligned}$$

which can be multiplied by

$$\Pr\left(\bigcap_{k=1}^{t-1} \Omega_1(k) \mid \bigcap_{k=1}^{t-1} \Omega_2(k), \bigcap_{k=1}^{t-1} \Omega_4(k), \bigcap_{k=1}^{t-1} \Omega_5(k)\right).$$

The result is then multiplied by

$$\Pr\left(\bigcap_{k=1}^{t-1} \Omega_2(k) \middle| \bigcap_{k=1}^{t-2} \Omega_1(k), \bigcap_{k=1}^{t-1} \Omega_3(k)\right) = 1.$$

Then the following equality holds:

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k) \middle| \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^{t-2} \Omega_1(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1$$

which uses the fact that  $\bigcap_{k=1}^t \Omega_y(k) \cap \bigcap_{k=1}^{t-1} \Omega_y(k) = \bigcap_{k=1}^t \Omega_y(k)$  ( $y = 1, 2, 3, 4, 5$ ). So by repeating this until the iteration  $j_0$ , then the following equality holds:

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k) \middle| \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^{j_0} \Omega_1(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1 \quad (58)$$

When  $t \leq j_0$ , we know that  $\mathbf{w}^{I,S}_t = \mathbf{w}^{U,S}_t$  and  $M_1^S \geq M_1$ , which means that if  $\Omega_0(k)$  holds, then  $\Omega_1(k)$  holds when  $\mathbf{w}^{I,S}_t = \mathbf{w}^{U,S}_t$ , and thus

$$\Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k) \middle| \bigcap_{k=1}^{j_0} \Omega_0(k)\right) = 1.$$

Then according to Theorem 8, we know that:

$$\Pr\left(\bigcap_{k=1}^{j_0} \Omega_0(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) = 1.$$

By multiplying the above two formulas, we get:

$$\begin{aligned} & \Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k) \middle| \bigcap_{k=1}^{j_0} \Omega_0(k)\right) \cdot \Pr\left(\bigcap_{k=1}^{j_0} \Omega_0(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) \\ &= \Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k), \bigcap_{k=1}^{j_0} \Omega_0(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) = 1 \end{aligned}$$

Note that since the probability of two joint events is smaller than that of either of the events, the following inequality holds:

$$\Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k), \bigcap_{k=1}^{j_0} \Omega_0(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) \leq \Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) \leq 1.$$

So we know that:

$$\Pr\left(\bigcap_{k=1}^{j_0} \Omega_1(k) \middle| \bigcap_{k=1}^{j_0} \Omega_5(k)\right) = 1.$$

which can be multiplied by Formula (58) and thus the following equality holds:

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k) \middle| \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1 \quad (59)$$

Then we can compute the probability of the negation of the joint event  $(\bigcap_{k=1}^{t+1} \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k))$ :

$$\begin{aligned}
& \Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)}\right) \\
&= \Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)} \middle| \overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \cdot \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \\
&+ \Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)} \middle| \overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \cdot \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \\
&\leq \Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)} \middle| \overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) + \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \\
&= \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right).
\end{aligned}$$

The last two steps use the fact that

$$\Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)} \middle| \overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) = 0$$

and

$$\Pr\left(\overline{\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)} \middle| \overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) \leq 1.$$

By further using the property of the probability of the union of multiply events, the formula above is bounded as:

$$\begin{aligned}
\Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)}\right) &\leq \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k)} \cup \overline{\bigcap_{k=1}^t \Omega_5(k)} \cup \overline{\bigcap_{k=1}^t \Omega_3(k)}\right) \\
&\leq \Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k)}\right) + \Pr\left(\overline{\bigcap_{k=1}^t \Omega_5(k)}\right) + \Pr\left(\overline{\bigcap_{k=1}^t \Omega_3(k)}\right).
\end{aligned}$$

Then by using Theorem 6, Formula (53), Lemma 9 and taking the union between iteration 0 and  $t$ , we get:

$$\begin{aligned}
\Pr\left(\overline{\bigcap_{k=1}^t \Omega_3(k)}\right) &\leq t \times \Psi_2, \\
\Pr\left(\overline{\bigcap_{k=1}^t \Omega_4(k)}\right) &\leq t \Psi_3, \\
\Pr\left(\overline{\bigcap_{k=1}^t \Omega_5(k)}\right) &\leq t \times \Psi_1.
\end{aligned}$$

Then we can know that:

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)\right) \geq 1 - t(\Psi_2 + \Psi_3 + \Psi_1)$$

and thus

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k)\right) \geq \Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k)\right) \geq 1 - t(\Psi_2 + \Psi_3 + \Psi_1).$$

This finishes the proof.

Similarly, from Formula (59), we know that for all  $T$  iterations:

$$\Pr\left(\bigcap_{k=1}^t \Omega_1(k), \bigcap_{k=1}^t \Omega_2(k), \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k) \mid \bigcap_{k=1}^t \Omega_4(k), \bigcap_{k=1}^t \Omega_5(k), \bigcap_{k=1}^t \Omega_3(k)\right) = 1. \quad (60)$$

Through the same argument, we know that:

$$\Pr\left(\bigcap_{k=1}^T \Omega_1(k), \bigcap_{k=1}^T \Omega_2(k), \bigcap_{k=1}^T \Omega_4(k), \bigcap_{k=1}^T \Omega_5(k)\right) \geq 1 - T(\Psi_2 + \Psi_3 + \Psi_1).$$

□

### A.3.7 Proof of Theorem 11

The proof is the same as the proof of Theorem 4 except that  $\mathbf{w}$ ,  $\mathbf{w}^I$ ,  $n$ ,  $r$ ,  $\xi_{j_1, j_m}$ ,  $\mathbf{H}$ ,  $\mathbf{B}$  need to be replaced by  $\mathbf{w}^S$ ,  $\mathbf{w}^{I,S}$ ,  $B$ ,  $r$ ,  $\xi_{j_1, j_m}^S$ ,  $\mathbf{H}^S$ ,  $\mathbf{B}^S$  and the main theorems that the proof depends on will be replaced by Theorem 9 and Corollary 2. We will show some key steps below.

First of all, according to the proofs of Theorem 10, we know that the following inequalities hold with probability higher than  $1 - T(\Psi_2 + \Psi_3 + \Psi_1)$ :

$$\begin{aligned} \|\mathbf{w}^S_k - \mathbf{w}^{I,S}_k\| &\leq \frac{1}{\frac{1}{2} - \frac{r}{n} - \frac{1}{B^{1/4}}} M_1 \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right); \\ \|\mathbf{H}^S_k - \mathbf{B}^S_{j_m}\| &\leq \xi_{j_1, j_m}^S; \\ \xi_{j_0+xT_0, j_0+(x+m-1)T_0}^S &\leq (1 - \eta\mu)^{j_0+xT_0} A d_{0, mT_0-1} + A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4} \\ &\leq \xi_{j_0, j_0+(m-1)T_0} + AM_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}; \\ \frac{\Delta B_k}{B} &\leq \frac{r}{n} + \frac{1}{B^{1/4}}. \end{aligned}$$

Then by subtracting  $\mathbf{w}^{I,S}_t$  by  $\mathbf{w}^{U,S}_t$  and following the arguments from Formula (32) to (34), the following inequality holds for  $\|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\|$  with probability higher than  $1 - T \times (\Psi_2 + \Psi_1 + \Psi_3)$ :

$$\begin{aligned} &\|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| \\ &\leq (1 - \eta\mu + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S) \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| \\ &+ \frac{\eta c_0}{2} \|\mathbf{w}^{U,S}_t - \mathbf{w}^{I,S}_t\| \|\mathbf{w}^{U,S}_t - \mathbf{w}^S_t\| + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S \|\mathbf{w}^{U,S}_t - \mathbf{w}^S_t\| \\ &\leq \left( 1 - \eta\mu + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S + \frac{c_0 M_1 \eta}{2} \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right) \right) \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S M_1 \left( \frac{r}{n} + \frac{1}{B^{1/4}} \right). \end{aligned}$$

By using the fact that  $\frac{\Delta B_t}{B} \leq \frac{r}{n} + \frac{1}{B^{1/4}}$  and  $\xi_{j_1, j_m} \leq \xi_{j_0, j_0+(m-1)T_0} + A \times M_1 \left( \frac{(\log(p+1))^2}{B} \right)^{1/4}$ , the formula



above can be bounded as:

$$\begin{aligned}
& \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| \\
& \leq \left(1 - \eta\mu + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S + \frac{c_0 M_1 \eta}{2} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)\right) \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| + \frac{B\eta}{B - \Delta B_t} \xi_{j_1, j_m}^S M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right) \\
& \leq \left[1 - \eta\mu + \frac{\eta}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} (\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4})\right] \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| \\
& \quad + \frac{c_0 M_1 \eta}{2} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right) \|\mathbf{w}^{I,S}_t - \mathbf{w}^{U,S}_t\| + \frac{\eta}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} \xi_{j_0+xT_0, j_0+(x+m-1)T_0}^S M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right).
\end{aligned}$$

Since  $\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4} \leq \frac{\mu}{2}$  and  $B$  is a large mini-batch size, then

$$1 - \left(\eta\mu - \frac{\eta}{1 - \frac{r}{n} - \frac{1}{B^{1/4}}} (\xi_{j_0, j_0+(m-1)T_0}^S + A \times M_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}) - \frac{c_0 M_1 \eta}{2} \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)\right) < 1.$$

Then after explicitly using the definition of  $\xi_{j_1, j_m}^S$  and following the argument of equation (35) to (38), we get:

$$\begin{aligned}
& \|\mathbf{w}^{I}_{j_0+(y+m)T_0} - \mathbf{w}^{U}_{j_0+(y+m)T_0}\| \\
& \leq (1 - \eta C)^{yT_0} \|\mathbf{w}^{I}_{j_0+mT_0} - \mathbf{w}^{U}_{j_0+mT_0}\| \\
& \quad + \frac{M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)}{C \left(1 - \frac{r}{n} - \frac{1}{B^{1/4}}\right)} (1 - \eta C)^{yT_0} (1 - \eta\mu)^{j_0} d_{0, mT_0-1} \frac{1}{1 - \left(\frac{1-\eta\mu}{1-\eta C}\right)^{T_0}} \\
& \quad + \frac{1}{1 - (1 - \eta C)^{T_0}} \frac{M_1 \left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)}{C \left(1 - \frac{r}{n} - \frac{1}{B^{1/4}}\right)} \left(A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}\right)
\end{aligned} \tag{61}$$

when  $t \rightarrow \infty$  and thus  $y \rightarrow \infty$ ,  $(1 - \eta C)^{yT_0} \rightarrow 0$ . Also with large mini-batch value  $B$ ,  $A_1 \frac{r}{n} + A_2 \frac{1}{B^{1/4}} + AM_1 \left(\frac{(\log(p+1))^2}{B}\right)^{1/4}$  is a value of the same order as  $\frac{r}{n} + \frac{1}{B^{1/4}}$ . Thus

$$\|\mathbf{w}^{I}_{j_0+(y+m)T_0} - \mathbf{w}^{U}_{j_0+(y+m)T_0}\| = o\left(\frac{r}{n} + \frac{1}{B^{1/4}}\right)$$

and

$$\|\mathbf{w}^{U,S}_t - \mathbf{w}^{I,S}_t\| \leq o\left(\frac{r}{n} + \frac{1}{B^{1/4}}\right).$$

## B Details on applications

### B.1 Privacy related data deletion

The notion of Approximate Data Deletion from the training dataset is proposed in Ginart et al. (2019):

**Definition 1.** A data deletion operation  $R_A$  is a  $\delta$ -deletion for algorithm  $A$  if, for all datasets  $D$  and for all measurable subset  $S$ , the following inequality holds:

$$Pr[A(D_{-i}) \in S | D_{-i}] \geq \delta Pr[R_A(D, A(D), i) \in S | D_{-i}],$$

where  $D$  is the full training dataset,  $D_{-i}$  is the remaining dataset after the  $i$ th sample is removed,  $A(D)$  and  $A(D_{-i})$  represent the model trained over  $D$  and  $D_{-i}$  respectively. Also  $R_A$  is an approximate model update algorithm, which updates the model after the sample  $i$  is removed.

This definition mimics the classical definition of differential privacy (Dwork et al., 2014):

**Definition 2.** A mechanism  $M$  is  $\epsilon$ -differentially private, where  $\epsilon \geq 0$ , if for all neighboring databases  $D_0$  and  $D_1$ , i.e., for databases differing in only one record, and for all sets  $S \in [M]$ , where  $[M]$  is the range of  $M$ , the following inequality holds:

$$\Pr[M(D_0) \in S] \leq e^\epsilon \Pr[M(D_1) \in S].$$

By borrowing the notations from Ginart et al. (2019), we define a version of approximate data deletion, which is slightly more strict than the one from Ginart et al. (2019):

**Definition 3.**  $R_A$  is an  $\epsilon$ -approximate deletion for  $A$  if for all  $D$  and measurable subset  $S \subset \mathcal{H}$ :

$$P(A(D_{-i}) \in S | D_{-i}) \leq e^\epsilon P(R_A(D, A(D), i) \in S | D_{-i})$$

and

$$P(R_A(D, A(D), i) \in S | D_{-i}) \leq e^\epsilon P(A(D_{-i}) \in S | D_{-i}).$$

To satisfy this definition for gradient descent, necessary randomness is added to the output of the BaseL and DeltaGrad. One simple way is the Laplace mechanism (Dwork et al., 2014), also following the idea from Chaudhuri and Monteleoni (2009) where noise following the Laplace distribution, i.e.

$$\text{Lap}(x | \frac{2}{n\epsilon\lambda}) = \frac{1}{2} \exp(-\frac{|x|}{\frac{2}{n\epsilon\lambda}}),$$

is added to the each coordinate of the output of the regularized logistic regression. Here  $p$  is the number of the parameters,  $\lambda$  is the regularization rate and  $\frac{2}{n\lambda}$  is the *sensitivity* of logistic regression (see Chaudhuri and Monteleoni (2009) for more details).

We can add even smaller noise to  $\mathbf{w}^*$ ,  $\mathbf{w}^{U^*}$  and  $\mathbf{w}^{I^*}$ , which follows the distribution  $\text{Lap}(\frac{\delta}{\epsilon})$  for each coordinate of  $\mathbf{w}^*$ ,  $\mathbf{w}^{U^*}$  and  $\mathbf{w}^{I^*}$  and is independent across different coordinates. Here  $\delta > \sqrt{p}\delta_0$  and

$$\delta_0 = \frac{1}{\eta(\frac{1}{2}\mu - \frac{r}{n-r}\mu - \frac{c_0 M_1 r}{2n})^2} \frac{M_1 r}{n-r} (A \frac{1}{\frac{1}{2} - \frac{r}{n}} M_1 \frac{r}{n})$$

(which is an upper bound on  $\|\mathbf{w}^{U^*} - \mathbf{w}^{I^*}\|$ ), such that the randomized DeltaGrad preserves  $\epsilon$ -approximate deletion.

*Proof.* We denote the model parameters after adding the random noise over  $\mathbf{w}^R$ ,  $\mathbf{w}^{U,R}$  and  $\mathbf{w}^{I,R}$ , and  $\mathbf{v}_i$  as the value of  $\mathbf{v}$  in the  $i_{th}$  coordinate. We have:

$$\mathbf{w}^* - \mathbf{w}^{R^*}, \mathbf{w}^{U^*} - \mathbf{w}^{U,R^*}, \mathbf{w}^{I^*} - \mathbf{w}^{I,R^*} \sim \text{Lap}(\frac{\delta}{\epsilon})$$

Given an arbitrary vector  $\mathbf{z} = [z_1, z_2, \dots, z_p]$ , the probability density ratio between  $\text{Pdf}(\mathbf{w}^{U,R^*} = \mathbf{z})$  and  $\text{Pdf}(\mathbf{w}^{I,R^*} = \mathbf{z})$  can be calculated as

$$\begin{aligned} \frac{\text{Pdf}(\mathbf{w}^{U,R^*} = \mathbf{z})}{\text{Pdf}(\mathbf{w}^{I,R^*} = \mathbf{z})} &= \frac{\prod_{i=1}^p \frac{\epsilon}{\delta} \exp(-\frac{\epsilon|\mathbf{z}_i - \mathbf{w}^{U^*}_i|}{\delta})}{\prod_{i=1}^p \frac{\epsilon}{\delta} \exp(-\frac{\epsilon|\mathbf{z}_i - \mathbf{w}^{I^*}_i|}{\delta})} \\ &= \prod_{i=1}^p \exp(\frac{\epsilon(|\mathbf{z}_i - \mathbf{w}^{U^*}_i| - |\mathbf{z}_i - \mathbf{w}^{I^*}_i|)}{\delta}) \\ &\leq \prod_{i=1}^p \exp(\frac{\epsilon(|\mathbf{w}^{I^*}_i - \mathbf{w}^{U^*}_i|)}{\delta}) \\ &= \exp(\frac{\epsilon(\|\mathbf{w}^{I^*} - \mathbf{w}^{U^*}\|_1)}{\delta}) \end{aligned}$$

Since

$$\|\mathbf{w}^{I^*} - \mathbf{w}^{U^*}\|_1 \leq \sqrt{p} \|\mathbf{w}^{I^*} - \mathbf{w}^{U^*}\|_2 = \sqrt{p} \|\mathbf{w}^{I^*} - \mathbf{w}^{U^*}\|$$

Then,

$$\begin{aligned} \frac{Pdf(\mathbf{w}^{U,R^*} = \mathbf{z})}{Pdf(\mathbf{w}^{I,R^*} = \mathbf{z})} &\leq \exp\left(\frac{\epsilon(\|\mathbf{w}^{I^*} - \mathbf{w}^{U^*}\|)}{\delta}\right) \\ &\leq \exp\left(\frac{\epsilon\sqrt{p}\delta_0}{\delta}\right) \leq \exp(\epsilon) \end{aligned}$$

Similarly, we can also prove  $\frac{Pdf(\mathbf{w}^{U,R^*} = \mathbf{z})}{Pdf(\mathbf{w}^{I,R^*} = \mathbf{z})} \geq \exp(\epsilon)$  by symmetry. □

## C Supplementary algorithm details

In Section 2, we only provided the details of DeltaGrad for deterministic gradient descent for the strongly convex and smooth objective functions in batch deletion/addition scenarios. In this section, we will provide more details on how to extend DeltaGrad to handle stochastic gradient descent, online deletion/addition scenarios and non-strongly convex, non-smooth objective functions.

### C.1 Extension of DeltaGrad for stochastic gradient descent

By using the notations from equations (5)-(7), we need to approximately or explicitly compute  $G_{B,S}$ , i.e. the average gradient for a mini-batch in the SGD version of DeltaGrad, instead of  $\nabla F$ , which is the average gradient for all samples. So by replacing  $\mathbf{w}_t$ ,  $\mathbf{w}^U_t$ ,  $\mathbf{w}^I_t$ ,  $\nabla F$ ,  $\mathbf{B}$  and  $\mathbf{H}$  with  $\mathbf{w}^S_t$ ,  $\mathbf{w}^{U,S}_t$ ,  $\mathbf{w}^{I,S}_t$ ,  $G_{B,S}$ ,  $\mathbf{B}^S$  and  $\mathbf{H}^S$  in Algorithm 1, we get the SGD version of DeltaGrad.

### C.2 Extension of DeltaGrad for online deletion/addition

In the online deletion/addition scenario, whenever the model parameters are updated after the deletion or addition of one sample, the history information should be also updated to reflect the changes. By assuming that only one sample is deleted or added each time, the online deletion/addition version of DeltaGrad is provided in Algorithm 2 and the differences relative to Algorithm 1 are highlighted.

Since the history information needs to be updated every time when new deletion or addition requests arrive, we need to do some more analysis on the error bound, which is still pretty close to the analysis in Section A.

In what follows, the analysis will be conducted on gradient descent with online deletion. Other similar scenarios, e.g. stochastic gradient descent with online addition, will be left as the future work.

#### C.2.1 Convergence rate analysis for online gradient descent version of DeltaGrad

##### Additional notes on setup, preliminaries

Let us still denote the model parameters for the original dataset at the  $t_{th}$  iteration by  $\mathbf{w}_t$ . During the model update phase for the  $k_{th}$  deletion request at the  $t_{th}$  iteration, the model parameters updated by BaseL and DeltaGrad are denoted by  $\mathbf{w}^U_t(k)$  and  $\mathbf{w}^I_t(k)$  respectively where  $\mathbf{w}^U_t(0) = \mathbf{w}^I_t(0) = \mathbf{w}_t$ . We also assume that the total number of removed samples in all deletion requests,  $r$ , is still far smaller than the total number of samples,  $n$ .

Also suppose that the indices of the removed samples are  $\{i_1, i_2, \dots, i_r\}$ , which are removed at the  $1_{st}$ ,  $2_{nd}$ ,  $3_{rd}$ ,  $\dots$ ,  $r_{th}$  deletion request. This also means that the cumulative number of samples up to the  $k_{th}$  deletion request ( $k \leq r$ ) is  $n - k$  for all  $1 \leq k \leq r$  and thus the objective function at the  $k_{th}$  iteration will be:

$$F^k(\mathbf{w}) = \frac{1}{n - k} \sum_{i \notin R_k} F_i(\mathbf{w}).$$

where  $R_k = \{i_1, i_2, \dots, i_k\}$ . Plus, at the  $k_{th}$  deletion request, we denote by  $\mathbf{H}_t^k$  the average Hessian matrix of  $F^k(\mathbf{w})$  evaluated between  $\mathbf{w}^I_t(k+1)$  and  $\mathbf{w}^I_t(k)$ :

$$\mathbf{H}_t^k = \frac{1}{n - k} \sum_{i \notin R_k} \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(k) + x(\mathbf{w}^I_t(k+1) - \mathbf{w}^I_t(k))) dx$$

---

**Algorithm 2:** DeltaGrad (online deletion/addition)

---

**Input** : The full training set  $(\mathbf{X}, \mathbf{Y})$ , model parameters cached during the training phase for the full training samples  $\{\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_t\}$  and corresponding gradients  $\{\nabla F(\mathbf{w}_0), \nabla F(\mathbf{w}_1), \dots, \nabla F(\mathbf{w}_t)\}$ , the index of the removed training sample or the added training sample  $i_r$ , period  $T_0$ , total iteration number  $T$ , history size  $m$ , warmup iteration number  $j_0$ , learning rate  $\eta$

**Output:** Updated model parameter  $\mathbf{w}^I_t$

- 1 Initialize  $\mathbf{w}^I_0 \leftarrow \mathbf{w}_0$
- 2 Initialize an array  $\Delta G = []$
- 3 Initialize an array  $\Delta W = []$
- 4 **for**  $t = 0; t < T; t++$  **do**
- 5     **if**  $[(t - j_0) \bmod T_0] == 0$  **or**  $t \leq j_0$  **then**
- 6         compute  $\nabla F(\mathbf{w}^I_t)$  exactly
- 7         compute  $\nabla F(\mathbf{w}^I_t) - \nabla F(\mathbf{w}_t)$  based on the cached gradient  $\nabla F(\mathbf{w}_t)$
- 8         set  $\Delta G[k] = \nabla F(\mathbf{w}^I_t) - \nabla F(\mathbf{w}_t)$
- 9         set  $\Delta W[k] = \mathbf{w}^I_t - \mathbf{w}_t$ , based on the cached parameters  $\mathbf{w}_t$
- 10          $k \leftarrow k + 1$
- 11         compute  $\mathbf{w}^I_{t+1}$  by using exact GD update (equation (1))
- 12          $\mathbf{w}_t \leftarrow \mathbf{w}^I_t$
- 13          $\nabla F(\mathbf{w}_t) \leftarrow \nabla F(\mathbf{w}^I_t)$
- 14     **else**
- 15         Pass  $\Delta W[-m:], \Delta G[-m:]$ , the last  $m$  elements in  $\Delta W$  and  $\Delta G$ , which are from the  $j_1^{th}, j_2^{th}, \dots, j_m^{th}$  iterations where  $j_1 < j_2 < \dots < j_m$  depend on  $t$ ,  $\mathbf{v} = \mathbf{w}^I_t - \mathbf{w}_t$ , and the history size  $m$ , to the L-BFGFS Algorithm (See Supplement) to get the approximation of  $\mathbf{H}(\mathbf{w}_t)\mathbf{v}$ , i.e.,  $\mathbf{B}_{j_m}\mathbf{v}$
- 16         Approximate  $\nabla F(\mathbf{w}^I_t) = \nabla F(\mathbf{w}_t) + \mathbf{B}_{j_m}(\mathbf{w}^I_t - \mathbf{w}_t)$
- 17         Compute  $\mathbf{w}^I_{t+1}$  by using the "leave-1-out" gradient formula, based on the approximated  $\nabla F(\mathbf{w}^I_t)$
- 18          $\mathbf{w}_t \leftarrow \mathbf{w}^I_t$
- 19          $\nabla F(\mathbf{w}_t) \leftarrow \frac{\eta}{n-1}[n(\mathbf{B}_{j_m}(\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)) - \nabla F_{i_r}(\mathbf{w}_t)]$
- 20     **end**
- 21 **end**
- 22 **return**  $\mathbf{w}^I_t$

---

Specifically,

$$\mathbf{H}_t^0 = \frac{1}{n} \sum_{i=1}^n \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(0) + x(\mathbf{w}^I_t(1) - \mathbf{w}^I_t(0))) dx.$$

Also the model parameters and the approximate gradients evaluated by DeltaGrad at the  $r - 1_{st}$  deletion request are used at the  $r_{th}$  request, and are denoted by:

$$\{\mathbf{w}^I_0(r-1), \mathbf{w}^I_1(r-1), \dots, \mathbf{w}^I_t(r-1)\}$$

and

$$\{g^a(\mathbf{w}^I_0(r-1)), g^a(\mathbf{w}^I_1(r-1)), \dots, g^a(\mathbf{w}^I_t(r-1))\}.$$

Note that  $g^a(\mathbf{w}^I_t(k))$  ( $k \leq r$ ) is not necessarily equal to  $\nabla F$  due to the approximation brought by DeltaGrad. But due to the periodicity of DeltaGrad, at iteration  $0, 1, \dots, j_0$  and iteration  $j_0 + xT_0$  ( $x = 1, 2, \dots$ ), the gradients are explicitly evaluated, i.e.:

$$g^a(\mathbf{w}^I_t(k)) = \frac{1}{n-k} \sum_{i \notin R_k} \nabla F_i(\mathbf{w}^I_t(k))$$

for  $t = 0, 1, \dots, j_0$  or  $t = j_0 + xT_0$  ( $x \geq 1$ ) and all  $k \leq r$ .

Also, due to the periodicity, the sequence  $[\Delta g_{j_0}, \Delta g_{j_1}, \dots, \Delta g_{j_{m-1}}]$  used in approximating the Hessian matrix always uses the exact gradient information, which means that:

$$\Delta g_{j_q} = \frac{1}{n-k} \left[ \sum_{i \notin R_k} \nabla F_i(\mathbf{w}^I_{j_q}(k)) - \sum_{i \notin R_k} \nabla F_i(\mathbf{w}^I_{j_q}(k-1)) \right]$$

where  $q = 1, 2, \dots, m-1$ . So Lemma 6 on the bound on the eigenvalues of  $B_{j_q}$  holds for all  $q$  and  $k = 1, 2, \dots, r$ .

But for the iterations where the gradients are not explicitly evaluated, the calculation of  $g^a(\mathbf{w}^I_t(k))$  depends on the approximated Hessian matrix  $\mathbf{B}_{j_m}^{k-1}$  and the approximated gradients calculated at the  $t_{th}$  iteration at the  $k-1$ -st deletion request. So the update rule for  $g^a(\mathbf{w}^I_t(k))$  is:

$$g^a(\mathbf{w}^I_t(k)) = \frac{1}{n-k} \{(n-k+1)[\mathbf{B}_{j_m}^{k-1}(\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)) + g^a(\mathbf{w}^I_t(k-1))] - \nabla F_{i_k}(\mathbf{w}^I_t(k))\}. \quad (62)$$

Here the product  $\mathbf{B}_{j_m}^{k-1} \cdot (\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1))$  approximates

$$\frac{1}{n-k+1} \sum_{i \notin R_{k-1}} \nabla F_i(\mathbf{w}^I_t(k)) - \nabla F_i(\mathbf{w}^I_t(k-1))$$

and  $g^a(\mathbf{w}^I_t(k-1))$  approximates  $\frac{1}{n-k+1} \sum_{i \notin R_{k-1}} \nabla F_i(\mathbf{w}^I_t(k-1))$ .

Similarly, the online version of  $\Delta w$  (at the  $k_{th}$  iteration) becomes:

$$\Delta w_{j_q}(k) = \mathbf{w}^I_{j_q}(k) - \mathbf{w}^I_{j_q}(k-1)$$

where  $q = 1, 2, \dots, m-1$ .

Similarly, we use  $d_{j_a, j_b}(k)$  to denote the value of the upper bound  $d$  on the distance between the iterates at the  $k_{th}$  deletion request and use  $\mathbf{B}_{j_m}^{k-1}$  to denote the approximated Hessian matrix in the  $k_{th}$  deletion request, which approximated the Hessian matrix  $\mathbf{H}_t^{k-1}$ .

So the update rule for  $\mathbf{w}^I_t(k)$  becomes:

$$\mathbf{w}^I_{t+1}(k) = \begin{cases} \mathbf{w}^I_t(k) - \frac{\eta}{n-k} \sum_{i \notin R_k} \nabla F_i(\mathbf{w}^I_t(k)), & [(t-j_0) \bmod T_0 = 0] \text{ or } t \leq j_0 \\ \mathbf{w}^I_t(k) - \frac{\eta}{n-k} \{(n-k+1)[\mathbf{B}_{j_m}^{k-1}(\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)) + g^a(\mathbf{w}^I_t(k-1))] - \nabla F_{i_k}(\mathbf{w}^I_t(k))\}, & \text{else.} \end{cases} \quad (63)$$

### Proof preliminaries.

On each deletion request, the BaseL model parameters are retrained from scratch on the remaining samples. This implies that Theorem 2 still holds, if we replace  $\mathbf{w}^U_t$ ,  $\mathbf{w}_t$  and  $r$  with  $\mathbf{w}^U_t(k)$ ,  $\mathbf{w}^U_t(k-1)$  and 1 respectively:

**Theorem 12** (Bound between iterates deleting one datapoint).  $\|\mathbf{w}^U_t(r) - \mathbf{w}^U_t(r-1)\| \leq M_1 \frac{1}{n}$  where  $M_1 = \frac{2}{\mu} c_2$  is some positive constant that does not depend on  $t$ . Here  $\mu$  is the strong convexity constant, and  $c_2$  is the bound on the individual gradients.

By induction, we have:

$$\|\mathbf{w}^U_t(r) - \mathbf{w}_t\| = \|\mathbf{w}^U_t(r) - \mathbf{w}^U_t(0)\| \leq M_1 \frac{r}{n}. \quad (64)$$

Then let us do some analysis on  $d_{j_a, j_b}(k)$ . We use the notation  $M_1^r \frac{1}{n}$  for  $\frac{\frac{2M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} \frac{(2L+\mu)}{\mu}}$ , where  $M_1^r$  is a constant which does not depend on  $k$ .

**Lemma 11.** If  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq \frac{2M_1}{1 - \frac{k+1}{n} - \frac{2(k-1)}{n} \frac{(2L+\mu)}{\mu}}$  for all  $k \leq r$ , then  $d_{j_a, j_b}(r) \leq d_{j_a, j_b}(0) + 2r \cdot M_1^r \frac{1}{n}$  where  $M_1$  is defined in Theorem 12.

*Proof.* Recall that  $d_{j_a, j_b}(k) = \max(\|\mathbf{w}^I_y(k) - \mathbf{w}^I_z(k)\|)_{j_a < y < z < j_b}$ . Then for two arbitrary iterations  $y, z$ , let us bound  $\|\mathbf{w}^I_y(k) - \mathbf{w}^I_z(k)\|$  as below:

$$\begin{aligned} & \|\mathbf{w}^I_y(k) - \mathbf{w}^I_z(k)\| \\ &= \|\mathbf{w}^I_y(k) - \mathbf{w}^I_z(k) + \mathbf{w}^I_y(k-1) - \mathbf{w}^I_z(k-1) + \mathbf{w}^I_z(k-1) - \mathbf{w}^I_y(k-1)\| \\ &\leq \|\mathbf{w}^I_y(k) - \mathbf{w}^I_y(k-1)\| + \|\mathbf{w}^I_z(k) - \mathbf{w}^I_z(k-1)\| + \|\mathbf{w}^I_z(k-1) - \mathbf{w}^I_y(k-1)\|. \end{aligned}$$

Then by using the bound on  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\|$ , the above formula leads to:

$$\leq 2 \cdot \frac{\frac{2M_1}{n}}{1 - \frac{k+1}{n} - \frac{2(k-1)}{n} \left(\frac{2L+\mu}{\mu}\right)} + \|\mathbf{w}^I_z(k-1) - \mathbf{w}^I_y(k-1)\|.$$

By using that  $\frac{\frac{2M_1}{n}}{1 - \frac{k+1}{n} - \frac{2(k-1)}{n} \left(\frac{2L+\mu}{\mu}\right)} \leq \frac{\frac{2M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} \left(\frac{2L+\mu}{\mu}\right)}$  and applying it recursively for  $k = 1, 2, \dots, r$ , we have:

$$\|\mathbf{w}^I_y(r) - \mathbf{w}^I_z(r)\| \leq 2r \cdot \frac{\frac{2M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} \left(\frac{2L+\mu}{\mu}\right)} + \|\mathbf{w}^I_z(0) - \mathbf{w}^I_y(0)\|.$$

Then by using the definition of  $d_{j_a, j_b}(k)$ , the following inequality holds:

$$d_{j_a, j_b}(r) \leq d_{j_a, j_b + T_0 - 1}(0) + 2r \cdot \frac{\frac{2M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} \left(\frac{2L+\mu}{\mu}\right)}.$$

Recalling the definition of  $M_1^r \frac{1}{n}$ , this is exactly the required result.  $\square$

We also mention that, since  $\frac{\frac{2M_1}{n}}{1 - \frac{k+1}{n} - \frac{2(k-1)}{n} \left(\frac{2L+\mu}{\mu}\right)} \leq \frac{\frac{2M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} \left(\frac{2L+\mu}{\mu}\right)}$ , then  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n}$  for any  $k \leq r$ .

**Theorem 13.** Suppose that at the  $k_{th}$  deletion request,  $\|\mathbf{w}^I_{j_q}(k) - \mathbf{w}^I_{j_q}(k-1)\| \leq M_1^r \frac{1}{n}$ , where  $q = 1, 2, \dots, m$  and  $M_1 = \frac{2c_2}{\mu}$ . Let  $e = \frac{L(L+1)+K_2L}{\mu K_1}$  for the upper and lower bounds  $K_1, K_2$  on the eigenvalues of the quasi-Hessian from Lemma 6, and for the Lipschitz constant  $c_0$  of the Hessian. For  $1 \leq z+1 \leq y \leq m$  we have:

$$\|\mathbf{H}_{j_z}^{k-1} - \mathbf{H}_{j_y}^{k-1}\| \leq c_0 d_{j_z, j_y}(k-1) + c_0 M_1^r \frac{1}{n}$$

and

$$\|\Delta g_{j_z} - \mathbf{B}_{j_y}^{k-1} \Delta w_{j_z}\| \leq [(1+e)^{y-z-1} - 1] \cdot c_0 (d_{j_z, j_y} + M_1^r \frac{1}{n}) \cdot s_{j_1, j_m}(k-1)$$

where  $s_{j_1, j_m}(k-1) = \max(\|\Delta w_a(k-1)\|)_{a=j_1, j_2, \dots, j_m} = \max(\|\mathbf{w}^I_a(k-1) - \mathbf{w}^I_a(k-2)\|)_{a=j_1, j_2, \dots, j_m}$ . Recall that  $d$  is defined as the maximum gap between the steps of the algorithm for the iterations from  $j_z$  to  $j_y$ :

$$d_{j_z, j_y}(k-1) = \max(\|\mathbf{w}^I_a(k-1) - \mathbf{w}^I_b(k-1)\|)_{j_z \leq a \leq b \leq j_y}. \quad (65)$$

*Proof.* Let us bound the difference between the averaged Hessians  $\|\mathbf{H}_{j_z}^{k-1} - \mathbf{H}_{j_y}^{k-1}\|$ , where  $1 \leq z < y \leq m$ , using their definition, as well as using Assumption 4 on the Lipschitzness of the Hessian. First we can get the following equality:

$$\begin{aligned} & \|\mathbf{H}_{j_y}^{k-1} - \mathbf{H}_{j_z}^{k-1}\| \\ &= \left\| \int_0^1 [\mathbf{H}(\mathbf{w}^I_{j_y}(k-1) + x(\mathbf{w}^I_{j_y}(k) - \mathbf{w}^I_{j_y}(k-1)))] dx \right. \\ & \quad \left. - \int_0^1 [\mathbf{H}(\mathbf{w}^I_{j_z}(k-1) + x(\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1)))] dx \right\| \\ &= \left\| \int_0^1 [\mathbf{H}(\mathbf{w}^I_{j_y}(k-1) + x(\mathbf{w}^I_{j_y}(k) - \mathbf{w}^I_{j_y}(k-1)))] \right. \\ & \quad \left. - \mathbf{H}(\mathbf{w}^I_{j_z}(k-1) + x(\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1)))] dx \right\| \end{aligned} \quad (66)$$

Then we can bound this as:

$$\begin{aligned}
&\leq c_0 \int_0^1 \|\mathbf{w}^I_{j_y}(k-1) + x(\mathbf{w}^I_{j_y}(k) - \mathbf{w}^I_{j_y}(k-1)) \\
&\quad - [\mathbf{w}^I_{j_z}(k-1) + x(\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1))]\| dx \\
&\leq c_0 \|\mathbf{w}^I_{j_y}(k-1) - \mathbf{w}^I_{j_z}(k-1)\| \\
&\quad + \frac{c_0}{2} \|\mathbf{w}^I_{j_y}(k) - \mathbf{w}^I_{j_y}(k-1) - (\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1))\| \\
&\leq c_0 \|\mathbf{w}^I_{j_y}(k-1) - \mathbf{w}^I_{j_z}(k-1)\| \\
&\quad + \frac{c_0}{2} \|\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1)\| + \frac{c_0}{2} \|\mathbf{w}^I_{j_y}(k) - \mathbf{w}^I_{j_y}(k-1)\| \\
&\leq c_0 d_{j_y, j_z}(k-1) + c_0 M_1^r \frac{1}{n} \leq c_0 d_{j_1, j_m + T_0 - 1}(k-1) + c_0 M_1^r \frac{1}{n}.
\end{aligned}$$

On the last line, we have used the definition of  $d_{j_z, j_y}$ , and the assumption on the boundedness of  $\|\mathbf{w}^I_{j_z}(k) - \mathbf{w}^I_{j_z}(k-1)\|$ .

Then by following the rest of the proof of Theorem 1, we get:

$$\|\Delta g_{j_z} - \mathbf{B}_{j_y} \Delta w_{j_z}\| \leq [(1+e)^{y-z-1} - 1] \cdot c_0 (d_{j_z, j_y}(k-1) + M_1^r \frac{1}{n}) \cdot s_{j_1, j_m}(k-1).$$

□

Similarly, the online version of Corollary 1 also holds by following the same derivation as the proof of Corollary 1 (except that  $r$ ,  $\xi_{j_1, j_m}$  and  $d_{j_1, j_m + T_0 - 1}$  is replaced by 1,  $\xi_{j_1, j_m}(k-1)$  and  $d_{j_1, j_m + T_0 - 1}(k-1)$  respectively), i.e.:

**Corollary 3** (Approximation accuracy of quasi-Hessian to mean Hessian (online deletion)). *Suppose that at the  $k_{th}$  deletion request,  $\|\mathbf{w}^I_{j_s}(k) - \mathbf{w}^I_{j_s}(k-1)\| \leq M_1^r \frac{1}{n}$  and  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n}$  where  $s = 1, 2, \dots, m$ . Then for  $j_m \leq t \leq j_m + T_0 - 1$ ,*

$$\|\mathbf{H}_t^{k-1} - \mathbf{B}_{j_m}^{k-1}\| \leq \xi_{j_1, j_m}(k-1) := A d_{j_1, j_m + T_0 - 1}(k-1) + A M_1^r \frac{1}{n}. \quad (67)$$

Recall that  $A = \frac{c_0 \sqrt{m} [(1+e)^m - 1]}{c_1} + c_0$ , where  $c_0$  is the Lipschitz constant of the Hessian,  $c_1$  is the "strong independence" constant from Assumption 5, and  $d_{j_1, j_m + T_0 - 1}(k-1)$  is the maximal gap between the iterates of the GD algorithm on the full data from  $j_1$  to  $j_m + T_0 - 1$  after the  $k-1$ -st deletion.

Based on this, let us derive a bound on  $\|\nabla F_i(\mathbf{w}^I_t(r))\|$ ,  $\|g^a(\mathbf{w}^I_t(r)) - \frac{1}{n-r} \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^I_t(r))\|$  and  $\|g^a(\mathbf{w}^I_t(r))\|$ .

**Lemma 12.** *Suppose we are at an iteration  $t$  such that  $j_m \leq t \leq j_m + T_0 - 1$ . If the following inequality holds for all  $k < r$ :*

$$\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n},$$

then the following inequality holds for all  $i = 1, 2, \dots, n$ :

$$\|\nabla F_i(\mathbf{w}^I_t(r-1))\| \leq M_1^r \frac{1}{n} Lr + c_2.$$

*Proof.* By adding and subtracting  $\nabla F_i(\mathbf{w}^I_t(r-2))$  inside  $\|\nabla F_i(\mathbf{w}^I_t(r-1))\|$ , we get:

$$\begin{aligned}
&\|\nabla F_i(\mathbf{w}^I_t(r-1))\| \\
&= \|\nabla F_i(\mathbf{w}^I_t(r-1)) - \nabla F_i(\mathbf{w}^I_t(r-2)) + \nabla F_i(\mathbf{w}^I_t(r-2))\| \\
&\leq \|\nabla F_i(\mathbf{w}^I_t(r-1)) - \nabla F_i(\mathbf{w}^I_t(r-2))\| + \|\nabla F_i(\mathbf{w}^I_t(r-2))\|
\end{aligned}$$

The last inequality uses the triangle inequality. Then by using the Cauchy mean value theorem, the upper bound on the eigenvalue of the Hessian matrix (i.e. Assumption 2) and the bound on  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\|$ , the formula above is bounded as (recall  $\mathbf{H}_i$  is an integrated Hessian):

$$\begin{aligned} &= \|\mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r-2) - \mathbf{w}^I_t(r-1))) \cdot (\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2))\| + \|\nabla F_i(\mathbf{w}^I_t(r-2))\| \\ &\leq LM_1^r \frac{1}{n} + \|\nabla F_i(\mathbf{w}^I_t(r-2))\|. \end{aligned}$$

By using this recursively, we get:

$$\leq \sum_{k=1}^{r-1} M_1^r \frac{1}{n} L + \|\nabla F_i(\mathbf{w}^I_t(0))\| \leq M_1^r \frac{1}{n} Lr + c_2.$$

□

**Lemma 13.** *If at a given iteration  $t$  such that  $j_m \leq t \leq j_m + T_0 - 1$ , for all  $k < r$ , the following inequalities hold:*

$$\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n}$$

and

$$\xi_{j_1, j_m}(k-1) \leq \frac{\mu}{2},$$

then we have

$$\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\| \leq rM_1^r \frac{1}{n} \mu.$$

*Proof.* First of all,  $\frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1))$  can be rewritten as below by using the Cauchy mean-value theorem:

$$\begin{aligned} \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) &= \frac{1}{n-r+1} \left[ \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-1)) - \nabla F_{i_{r-1}}(\mathbf{w}^I_t(r-1)) \right] \\ &= \frac{1}{n-r+1} \{ (n-r+2) [\mathbf{H}_t^{r-2} \times (\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2))] \\ &\quad + \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - \nabla F_{i_{r-1}}(\mathbf{w}^I_t(r-1)) \}. \end{aligned}$$

By subtracting the above formula from equation (62), i.e., the update rule for the approximate gradient, the norm of the approximation error between true and approximate gradients is:

$$\begin{aligned} &\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\| \\ &= \frac{1}{n-r+1} \left\| (n-r+2) (\mathbf{H}_t^{r-2} - \mathbf{B}_{j_m}^{r-2}) \times (\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)) \right\| \\ &\quad + \left\| \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - (n-r+2) g^a(\mathbf{w}^I_t(r-2)) \right\| \end{aligned}$$

Then by using the triangle inequality, Corollary 3 on the approximation accuracy of the quasi-Hessian (where the bound is in terms of  $\xi$ ), and the bound on  $\|\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)\|$ , the formula above is bounded as:

$$\begin{aligned} &\leq \frac{n-r+2}{n-r+1} \|\mathbf{H}_t^{r-2} - \mathbf{B}_{j_m}^{r-2}\| \|\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)\| \\ &\quad + \frac{1}{n-r+1} \left\| \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - (n-r+2) g^a(\mathbf{w}^I_t(r-2)) \right\| \tag{68} \\ &\leq \frac{n-r+2}{n-r+1} \xi_{j_1, j_m}(r-2) M_1^r \frac{1}{n} + \frac{n-r+2}{n-r+1} \left\| \frac{1}{n-r+2} \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - g^a(\mathbf{w}^I_t(r-2)) \right\| \end{aligned}$$



By using that  $\xi_{j_1, j_m}(r-2) \leq \frac{\mu}{2}$ , the formula above is bounded as:

$$\leq \frac{n-r+2}{n-r+1} \frac{\mu}{2} (M_1^r \frac{1}{n}) + \frac{n-r+2}{n-r+1} \left\| \frac{1}{n-r+2} \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - g^a(\mathbf{w}^I_t(r-2)) \right\|$$

We can use this recursively. Note that  $\nabla F(\mathbf{w}^I_t(0)) = g^a(\mathbf{w}^I_t(0))$ . In the end, we get the following inequality:

$$\leq \sum_{k=1}^{r-1} \frac{n-k}{n-r} \frac{\mu}{2} (M_1^r \frac{1}{n}) \leq M_1^r \frac{1}{n} \frac{\mu}{2} \sum_{k=1}^{r-1} \frac{n-k}{n-r}$$

Also for  $r \ll n$ ,  $\frac{n-k}{n-r} \leq 2$  (in fact we assumed  $r/n \leq \delta$  for a sufficiently small  $\delta$ , so this holds). So we get the bound  $rM_1^r \frac{1}{n} \mu$ .  $\square$

Note that for  $\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\|$ , we get a tighter bound when  $t \rightarrow \infty$  by using equation (68), Lemma 11 (i.e.  $d_{j_a, j_b}(r) \leq d_{j_a, j_b}(0) + 2r \cdot M_1^r \frac{1}{n}$ ) and Lemma 8 without using  $\xi_{j_1, j_m}(r-1) \leq \frac{\mu}{2}$ , which starts by bounding  $\xi_{j_1, j_m}(k-1)$  where  $k \leq r$ ,  $j_1 = j_0 + xT_0$  and  $j_m = j_0 + (x+m-1)T_0$ :

$$\begin{aligned} \xi_{j_1, j_m}(k-1) &= Ad_{j_1, j_m + T_0 - 1}(k-1) + AM_1^r \frac{1}{n} \\ &\leq Ad_{j_1, j_m + T_0 - 1}(0) + 2(k-1)A \cdot M_1^r \frac{1}{n} + AM_1^r \frac{1}{n} \\ &\leq A(1 - \mu\eta)^{j_0 + xT_0} d_{0, mT_0 - 1}(0) + A(2k-1)M_1^r \frac{1}{n}, \end{aligned} \tag{69}$$

which can be plugged into Equation (68), i.e.:

$$\begin{aligned} &\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\| \\ &\leq \frac{n-r+2}{n-r+1} \xi_{j_1, j_m}(r-2) M_1^r \frac{1}{n} \\ &\quad + \frac{n-r+2}{n-r+1} \left\| \frac{1}{n-r+2} \sum_{i \notin R_{r-2}} \nabla F_i(\mathbf{w}^I_t(r-2)) - g^a(\mathbf{w}^I_t(r-2)) \right\| \\ &\leq \sum_{k=1}^{r-1} \frac{n-k+1}{n-k} \xi_{j_1, j_m}(k-1) M_1^r \frac{1}{n} \\ &\leq \sum_{k=1}^{r-1} \frac{n-k+1}{n-k} [A(1 - \mu\eta)^{j_0 + xT_0} d_{0, mT_0 - 1}(0) + A(2k-1)M_1^r \frac{1}{n}] \cdot M_1^r \frac{1}{n} \\ &\leq 2A(1 - \mu\eta)^{j_0 + xT_0} d_{0, mT_0 - 1}(0) r M_1^r \frac{1}{n} + 2A(r M_1^r \frac{1}{n})^2 \end{aligned} \tag{70}$$

The last step uses that  $\frac{n-k+1}{n-k} \leq 2$  and  $\sum_{k=1}^{r-1} (2k-1) < \sum_{k=1}^r (2k-1) = r^2$ . So when  $t \rightarrow \infty$  and thus  $x \rightarrow \infty$ ,  $\left\| \frac{1}{n-r} \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^I_t(r)) - g^a(\mathbf{w}^I_t(r)) \right\| = o(\frac{r}{n})$ .

Then based on Lemma 12 and 13, the bound on  $\|g^a(\mathbf{w}^I_t(r))\|$  becomes:

$$\begin{aligned} &\|g^a(\mathbf{w}^I_t(r-1))\| \\ &= \|g^a(\mathbf{w}^I_t(r-1)) - \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) + \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1))\| \\ &\leq \|g^a(\mathbf{w}^I_t(r-1)) - \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1))\| + \left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) \right\| \\ &= rM_1^r \frac{1}{n} \mu + M_1^r \frac{1}{n} Lr + c_2 = (r\mu + Lr)M_1^r \frac{1}{n} + c_2 \end{aligned} \tag{71}$$

## Main results

**Theorem 14** (Bound between iterates on full data and incrementally updated ones (online deletions)). *Suppose that for any  $k < r$ ,  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n}$ . At the  $r$ th deletion request, consider an iteration  $t$  indexed with  $j_m$  for which  $j_m \leq t < j_m + T_0 - 1$ , and suppose that we are at the  $x$ -th iteration of full gradient updates, so  $j_1 = j_0 + xT_0$ ,  $j_m = j_0 + (m-1+x)T_0$ . Suppose that we have the bounds  $\|\mathbf{H}_t^{r-1} - \mathbf{B}_{j_m}^{r-1}\| \leq \xi_{j_1, j_m}(r-1) = Ad_{j_1, j_m + T_0 - 1}(r-1) + A(M_1^r \frac{1}{n})$  (where we recalled the definition of  $\xi$ ) and*

$$\xi_{j_1, j_m}(r-1) = Ad_{j_1, j_m + T_0 - 1}(r-1) + A(M_1^r \frac{1}{n}) \leq \frac{\mu}{2}$$

for all iterations  $x$ . Then

$$\|\mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1)\| \leq M_1^r \frac{1}{n}.$$

Recall that  $c_0$  is the Lipschitz constant of the Hessian,  $M_1$  and  $A$  are defined in Theorem 12 and Corollary 3 respectively, and do not depend on  $t$ ,

Then by using the same derivation as the proof of Theorem 4, we get the following results at the  $r$ th deletion request.

**Theorem 15** (Bound between iterates on full data and incrementally updated ones (all iterations, online deletion)). *At the deletion request  $r$ , if for all  $k < r$ ,  $\|\mathbf{w}^I_t(k) - \mathbf{w}^I_t(k-1)\| \leq M_1^r \frac{1}{n}$  holds, then for any  $j_m < t < j_m + T_0 - 1$ ,*

$$\|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)\| \leq M_1^r \frac{1}{n}$$

and

$$\|\mathbf{H}_t^{r-1} - \mathbf{B}_{j_m}^{r-1}\| \leq \xi_{j_1, j_m}(r-1) := Ad_{j_1, j_m + T_0 - 1}(r-1) + AM_1^r \frac{1}{n}$$

and

$$\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\| \leq rM_1^r \frac{1}{n} \mu$$

hold

Then by induction (the base case is similar to Theorem 4), we know that the following theorem holds for all iterations  $t$ :

**Theorem 16** (Bound between iterates on full data and incrementally updated ones (all iterations, all deletion requests, online deletion)). *At the  $r$ th deletion request, for any  $j_m < t < j_m + T_0 - 1$ ,*

$$\|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)\| \leq M_1^r \frac{1}{n}$$

and

$$\|\mathbf{H}_t^{r-1} - \mathbf{B}_{j_m}^{r-1}\| \leq \xi_{j_1, j_m}(r-1) := Ad_{j_1, j_m + T_0 - 1}(r-1) + AM_1^r \frac{1}{n}$$

and

$$\left\| \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) - g^a(\mathbf{w}^I_t(r-1)) \right\| \leq rM_1^r \frac{1}{n} \mu$$

hold

Then by induction (from the  $r$ th deletion request to the  $1_{st}$  deletion request), the following inequality holds:

$$\|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(0)\| = \|\mathbf{w}^I_t(r) - \mathbf{w}_t\| \leq r \cdot M_1^r \frac{1}{n}$$

Then by using equation (64), the following inequality holds:

$$\begin{aligned}
& \|\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r-1)\| = \|\mathbf{w}^U_t(r) - \mathbf{w}_t + \mathbf{w}_t - \mathbf{w}^I_t(r-1)\| \\
& \leq \|\mathbf{w}^U_t(r) - \mathbf{w}_t\| + \|\mathbf{w}_t - \mathbf{w}^I_t(r-1)\| \\
& \leq M_1 \frac{r}{n} + (r-1) \cdot M_1^r \frac{1}{n} := M_2 \frac{r}{n}
\end{aligned} \tag{72}$$

where  $M_2$  is a constant which does not depend on  $t$  or  $k$ .

In the end, we get a similar result for the bound on  $\|\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r)\|$ :

**Theorem 17** (Convergence rate of DeltaGrad (online deletion)). *At the  $r_{th}$  deletion request, for all iterations  $t$ , the result  $\mathbf{w}^I_t(r)$  of DeltaGrad, Algorithm 2, approximates the correct iteration values  $\mathbf{w}^U_t(r)$  at the rate*

$$\|\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r)\| = o\left(\frac{r}{n}\right).$$

So  $\|\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r)\|$  is of a lower order than  $\frac{r}{n}$ .

#### The proof of Theorem 14

*Proof.* Note that the approximated update rules for  $\mathbf{w}^I_t$  at the  $r_{th}$  and the  $(r-1)_{st}$  deletion request are:

$$\begin{aligned}
\mathbf{w}^I_{t+1}(r) &= \mathbf{w}^I_t(r) - \frac{\eta}{n-r} \{(n-r+1)[\mathbf{B}_{j_m}^{r-1}(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&+ g^a(\mathbf{w}^I_t(r-1))] - \nabla F_{i_r}(\mathbf{w}^I_t(r))\}
\end{aligned} \tag{73}$$

and

$$\begin{aligned}
\mathbf{w}^I_{t+1}(r-1) &= \mathbf{w}^I_t(r-1) - \frac{\eta}{n-r+1} \{(n-r+2)[\mathbf{B}_{j_m}^{r-2}(\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)) \\
&+ g^a(\mathbf{w}^I_t(r-2))] - \nabla F_{i_{r-1}}(\mathbf{w}^I_t(r-1))\}.
\end{aligned} \tag{74}$$

Note that since

$$\begin{aligned}
g^a(\mathbf{w}^I_t(r-1)) &= \frac{1}{n-r+1} \{(n-r+2)[\mathbf{B}_{j_m}^{r-2}(\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)) \\
&+ g^a(\mathbf{w}^I_t(r-2))] - \nabla F_{i_{r-1}}(\mathbf{w}^I_t(r-1))\},
\end{aligned}$$

then equation (74) can be rewritten as:

$$\begin{aligned}
\mathbf{w}^I_{t+1}(r-1) &= \mathbf{w}^I_t(r-1) - \frac{\eta}{n-r+1} \{(n-r+2)[\mathbf{B}_{j_m}^{r-2}(\mathbf{w}^I_t(r-1) - \mathbf{w}^I_t(r-2)) \\
&+ g^a(\mathbf{w}^I_t(r-2))] - \nabla F_{i_{r-1}}(\mathbf{w}^I_t(r-1))\} \\
&= \mathbf{w}^I_t(r-1) - \eta g^a(\mathbf{w}^I_t(r-1)).
\end{aligned} \tag{75}$$

Then by subtracting equation (74) from equation (75), the result becomes:

$$\begin{aligned}
& \mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1) \\
&= (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) - \frac{\eta}{n-r} \{(n-r+1)[\mathbf{B}_{j_m}^{r-1}(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&+ g^a(\mathbf{w}^I_t(r-1))] - \nabla F_{i_r}(\mathbf{w}^I_t(r))\} + \eta g^a(\mathbf{w}^I_t(r-1)).
\end{aligned}$$

Then by adding and subtracting  $\mathbf{H}_t^{r-1}$  and  $\frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1))$  in the formula above and rearranging the result properly, it becomes:

$$\begin{aligned}
& \mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1) \\
&= (\mathbf{I} - \eta \frac{n-r+1}{n-r} (\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})) (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&\quad - \frac{\eta}{n-r} \{ (n-r+1) [\mathbf{H}_t^{r-1} (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1))] \\
&\quad + \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1)) - \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1)) \\
&\quad + g^a(\mathbf{w}^I_t(r-1))] - \nabla F_{i_r}(\mathbf{w}^I_t(r)) \} + \eta g^a(\mathbf{w}^I_t(r-1)).
\end{aligned} \tag{76}$$

Then by using the fact that

$$\begin{aligned}
& \mathbf{H}_t^{r-1} (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) + \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1)) \\
&= \frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r))
\end{aligned}$$

and

$$\left( \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r)) \right) - \nabla F_{i_r}(\mathbf{w}^I_t(r)) = \sum_{i \notin R_r} \nabla F(\mathbf{w}^I_t(r)),$$

Equation (76) becomes:

$$\begin{aligned}
& \mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1) \\
&= (\mathbf{I} - \eta \frac{n-r+1}{n-r} (\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})) (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&\quad - \frac{\eta}{n-r} \left[ \sum_{i \notin R_r} \nabla F(\mathbf{w}^I_t(r)) - \sum_{i \notin R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1)) \right. \\
&\quad \left. + (n-r+1) g^a(\mathbf{w}^I_t(r-1)) \right] + \eta g^a(\mathbf{w}^I_t(r-1)).
\end{aligned} \tag{77}$$

Also note that by using the Cauchy mean-value theorem, the following equation holds:

$$\begin{aligned}
& \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^I_t(r)) - \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1)) \\
&= \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^I_t(r)) - \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^I_t(r-1)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)) \\
&= \left[ \sum_{i \notin R_r} \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1))) dx \right] (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)),
\end{aligned}$$

which can be plugged into equation (77), i.e.:

$$\begin{aligned}
& \mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1) \\
&= (\mathbf{I} - \eta \frac{n-r+1}{n-r} (\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})) (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&\quad - \frac{\eta}{n-r} \left\{ \left[ \sum_{i \notin R_r} \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1))) dx \right] \right. \\
&\quad \cdot (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)) \\
&\quad \left. + (n-r+1) g^a(\mathbf{w}^I_t(r-1)) \right\} + \eta g^a(\mathbf{w}^I_t(r-1)),
\end{aligned} \tag{78}$$

which can be rearranged as:

$$\begin{aligned}
& \mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1) \\
&= (\mathbf{I} - \eta \frac{n-r+1}{n-r} (\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})) (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\
&\quad - \frac{\eta}{n-r} \left\{ \left[ \sum_{i \notin R_r} \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1))) dx \right] \right. \\
&\quad \cdot (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)) \left. \right\} - \frac{\eta}{n-r} g^a(\mathbf{w}^I_t(r-1)).
\end{aligned} \tag{79}$$

Then by taking the matrix norm on both sides of equation (79) and using that  $\|\mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)))\| \geq \mu$  and  $\|\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1}\| \leq \xi_{j_1, j_m}(r-1)$ , equation (79) can be bounded as:

$$\begin{aligned}
& \|\mathbf{w}^I_{t+1}(r) - \mathbf{w}^I_{t+1}(r-1)\| \\
&\leq (1 - \eta\mu) \|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)\| \\
&\quad + \frac{(n-r+1)\eta}{n-r} \xi_{j_1, j_m}(r-1) \|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)\| \\
&\quad + \frac{\eta}{n-r} \|\nabla F_{i_r}(\mathbf{w}^I_t(r-1))\| + \left\| \frac{\eta}{n-r} g^a(\mathbf{w}^I_t(r-1)) \right\|.
\end{aligned}$$

Then by using Lemma 12 and equation (71), the formula above becomes:

$$\begin{aligned}
& \leq (1 - \eta\mu + \frac{(n-r+1)\eta}{n-r} \xi_{j_1, j_m}(r)) \|\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)\| \\
& \quad + \frac{\eta}{n-r} (M_1^r \frac{1}{n} L(r-1) + c_2) + \frac{\eta}{n-r} (M_1^r \frac{1}{n} (r-1)\mu + M_1^r \frac{1}{n} L(r-1) + c_2).
\end{aligned}$$

By using the bound on  $\xi_{j_1, j_m}(r)$  and applying the above formula recursively across all iterations, the formula above becomes:

$$\begin{aligned}
& \leq \frac{1}{\eta\mu - \frac{\eta(n-r+1)}{n-r} \frac{\mu}{2}} \left( \frac{\eta}{n-r} (M_1^r \frac{1}{n} L(r-1) + c_2) \right. \\
& \quad \left. + \frac{\eta}{n-r} (M_1^r \frac{1}{n} L(r-1) + M_1^r \frac{1}{n} \mu (r-1) + c_2) \right) \\
& = \frac{2}{(n-r-1)\mu} \left( (M_1^r \frac{1}{n} (2L(r-1) + (r-1)\mu) + 2c_2) \right).
\end{aligned}$$

Then by using that  $M_1 = \frac{2c_2}{\mu}$  and  $M_1^r \frac{1}{n} = \frac{2M_1}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} (\frac{2L+\mu}{\mu})}$ , the formula above can be rewritten as:

$$\begin{aligned}
& = \frac{2}{(n-r-1)\mu} \frac{\frac{2M_1(r-1)}{n} (2L+\mu) + \mu M_1 (1 - \frac{r+1}{n} - \frac{2(r-1)}{n} (\frac{2L+\mu}{\mu}))}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} (\frac{2L+\mu}{\mu})} \\
& = \frac{2 \frac{M_1}{n}}{1 - \frac{r+1}{n} - \frac{2(r-1)}{n} (\frac{2L+\mu}{\mu})} = M_1^r \frac{1}{n}.
\end{aligned}$$

This finishes the proof. □

### The proof of Theorem 17

*Proof.* Recall that the update rule for  $\mathbf{w}^U_t(r)$  is:

$$\mathbf{w}^U_{t+1}(r) = \mathbf{w}^U_t(r) - \eta \frac{1}{n-r} \sum_{i \notin R_r} \nabla F(\mathbf{w}^U_t(r))$$

and the update rule for  $\mathbf{w}^I_t(r)$  is (where the gradients are explicitly evaluated):

$$\mathbf{w}^I_{t+1}(r) = \mathbf{w}^I_t(r) - \frac{\eta}{n-r} [(n-r+1)(\mathbf{B}_{j_m}^{r-1}(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) + g^a(\mathbf{w}^I_t(r-1))) - \nabla F_{i_r}(\mathbf{w}^I_t(r))].$$

Then by subtracting  $\mathbf{w}^I_{t+1}(r)$  from  $\mathbf{w}^U_{t+1}(r)$ , we get:

$$\begin{aligned} & \|\mathbf{w}^I_{t+1}(r) - \mathbf{w}^U_{t+1}(r)\| \\ &= \|\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r) - \frac{\eta}{n-r} \{(n-r+1)[\mathbf{B}_{j_m}^{r-1}(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\ &+ g^a(\mathbf{w}^I_t(r-1))] - \nabla F_{i_r}(\mathbf{w}^I_t(r))\} + \frac{\eta}{n-r} \sum_{i \notin R_r} \nabla F(\mathbf{w}^U_t(r))\|. \end{aligned}$$

Then by bringing in  $\mathbf{H}_t^{r-1}$  and  $\frac{1}{n-r+1} \sum_{i \in R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1))$  into the formula above, we get:

$$\begin{aligned} &= \|\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r) - \frac{(n-r+1)\eta}{n-r} [(\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\ &+ \mathbf{H}_t^{r-1} \times (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) + g^a(\mathbf{w}^I_t(r-1)) \\ &- \frac{1}{n-r+1} \sum_{i \in R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1)) + \frac{1}{n-r+1} \sum_{i \in R_{r-1}} \nabla F(\mathbf{w}^I_t(r-1))] \\ &+ \frac{\eta}{n-r} [\nabla F_{i_r}(\mathbf{w}^I_t(r)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)) + \nabla F_{i_r}(\mathbf{w}^I_t(r-1))] + \frac{\eta}{n-r} \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^U_t(r))\|. \end{aligned}$$

Then by using the triangle inequality and the result from equation (70), the formula above can be bounded as:

$$\begin{aligned} &\leq \|\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r) - \frac{(n-r+1)\eta}{n-r} [(\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) \\ &+ \mathbf{H}_t^{r-1} \times (\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1)) + \frac{1}{n-r+1} \sum_{i \in R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1))\| \\ &+ \frac{\eta}{n-r} [\nabla F_{i_r}(\mathbf{w}^I_t(r)) - \nabla F_{i_r}(\mathbf{w}^I_t(r-1)) + \nabla F_{i_r}(\mathbf{w}^I_t(r-1))] \\ &+ \frac{\eta}{n-r} \sum_{i \notin R_r} \nabla F_i(\mathbf{w}^U_t(r))\| + 2A(1 - \mu\eta)^{j_0+xT_0} d_{0,(m-1)T_0}(0) r M_1^r \frac{1}{n} + 2A(r M_1^r \frac{1}{n})^2. \end{aligned}$$

Note that the first matrix norm in this formula is the same as equation (33) by replacing  $n, r, \mathbf{w}^I_t, \mathbf{w}^U_t, \mathbf{w}_t, \mathbf{B}_{j_m}, \mathbf{H}_t$  and  $\nabla F(\mathbf{w}_t)$  with  $n-r+1, 1, \mathbf{w}^I_t(r), \mathbf{w}^U_t(r), \mathbf{w}^I_t(r-1), \mathbf{B}_{j_m}^{r-1}, \mathbf{H}_t^{r-1}$  and  $\frac{1}{n-r+1} \sum_{i \notin R_{r-1}} \nabla F_i(\mathbf{w}^I_t(r-1))$  reps.. So by following the same derivation, the formula above can be bounded as:

$$\begin{aligned} &\leq \|(\mathbf{I} - \frac{\eta}{n-r} \sum_{i \notin R_r} \mathbf{H}_{t,i}^{r-1})(\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r))\| \\ &+ \|\frac{(n-r+1)\eta}{n-r} [(\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})(\mathbf{w}^I_t(r) - \mathbf{w}^U_t(r))]\| \\ &+ \|\frac{\eta}{n-r} [\sum_{i \notin R_r} \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r-1))) dx \\ &- \int_0^1 \mathbf{H}_i(\mathbf{w}^I_t(r-1) + x(\mathbf{w}^I_t(r) - \mathbf{w}^I_t(r-1))) dx](\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r-1))\| \\ &+ \|\frac{(n-r+1)\eta}{n-r} [(\mathbf{B}_{j_m}^{r-1} - \mathbf{H}_t^{r-1})(\mathbf{w}^U_t(r) - \mathbf{w}^I_t(r-1))]\| \\ &+ 2A(1 - \mu\eta)^{j_0+xT_0} d_{0,(m-1)T_0}(0) r M_1^r \frac{1}{n} + 2A(r M_1^r \frac{1}{n})^2. \end{aligned}$$

Then by using the following facts:

1.  $\|\mathbf{I} - \eta \mathbf{H}_{t,i}^{r-1}\| \leq 1 - \eta\mu$ ;
2. from Theorem 16 on the approximation accuracy of the quasi-Hessian to mean Hessian, we have the error bound  $\|\mathbf{H}_t^{r-1} - \mathbf{B}_{j_m}^{r-1}\| \leq \xi_{j_1, j_m}(r-1)$ ;
3. we bound the difference of integrated Hessians using the strategy from Equation (20);
4. from Equation (72), we have the error bound  $\|\mathbf{w}_t^U(r) - \mathbf{w}_t^I(r-1)\| \leq M_2 \frac{r}{n}$  (and this requires no additional assumptions),

the expression can be bounded as follows:

$$\begin{aligned} &\leq (1 - \eta\mu + \frac{(n-r+1)\eta}{n-r} \xi_{j_0, j_0+(m-1)T_0}(r-1) + \frac{c_0 M_2 r \eta}{2n}) \|\mathbf{w}_t^I - \mathbf{w}_t^U\| \\ &+ \frac{M_2(n-r+1)r\eta}{n(n-r)} \xi_{j_1, j_m}(r-1) + 2A(1-\mu\eta)^{j_0+xT_0} d_{0, (m-1)T_0}(0) r M_1^r \frac{1}{n} \\ &+ 2A(r M_1^r \frac{1}{n})^2, \end{aligned}$$

which is very similar to equation (36) (except the difference in the coefficient). So by following the derivation after equation (36), we know that:

$$\|\mathbf{w}_t^I(r) - \mathbf{w}_t^U(r)\| = o\left(\frac{r}{n}\right)$$

when  $t \rightarrow \infty$ .

□

### C.3 Extension of DeltaGrad for non-strongly convex, non-smooth objective functions

For the original version of the L-BFGS algorithm, strong convexity is essential to make the secant condition hold. In this subsection, we present our extension of DeltaGrad to non-strongly convex, non-smooth objectives.

To deal with non-strongly convex objectives, we assume that convexity holds in some local regions. When constructing the arrays  $\Delta G$  and  $\Delta W$ , only the model parameters and their gradients where local convexity holds are used.

For local non-smoothness, we found that even a small distance between  $\mathbf{w}_t$  and  $\mathbf{w}_t^I$  can make the estimated gradient  $\nabla F(\mathbf{w}_t^I)$  drift far away from  $\nabla F(\mathbf{w}_t)$ . To deal with this, we explicitly check if the norm of  $\mathbf{B}_{j_m}(\mathbf{w}_t - \mathbf{w}_t^I)$  (which equals to  $\nabla F(\mathbf{w}_t^I) - \nabla F(\mathbf{w}_t)$ ) is larger than the norm of  $L(\mathbf{w}_t - \mathbf{w}_t^I)$  for a constant  $L$ . In our experiments,  $L$  is configured as 1. The details of the modifications above are highlighted in Algorithm 3.

## D Supplementary experiments

In this section, we present some supplementary experiments that could not be presented in the paper due to space limitations.

### D.1 Experiments with large deletion rate

In this experiment, instead of deleting at most 1% of training samples each time as we did in Section 4 in the main paper, we vary the deletion rate from 0 to up to 20% on MNIST dataset and still compare the performance between DeltaGrad (with  $T_0$  as 5 and  $j_0$  as 10) and BaseL. All other hyper-parameters such as the learning rate and mini-batch size remain the same as in Section 4 in the main paper.

The experimental results in Figure 1 show that even with the largest deletion rate, i.e. 20%, DeltaGrad can still be 1.67x faster than BaseL (2.27s VS 1.53s) and the error bound between their resulting model

---

**Algorithm 3:** DeltaGrad (general models)

---

**Input** : The full training set  $(\mathbf{X}, \mathbf{Y})$ , model parameters cached during the training phase for the full training samples  $\{\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_t\}$  and corresponding gradients  $\{\nabla F(\mathbf{w}_0), \nabla F(\mathbf{w}_1), \dots, \nabla F(\mathbf{w}_t)\}$ , the removed training sample or the added training sample  $R$ , period  $T_0$ , total iteration number  $T$ , history size  $m$ , warmup iteration number  $j_0$ , learning rate  $\eta$

**Output:** Updated model parameter  $\mathbf{w}^I_t$

```
1 Initialize  $\mathbf{w}^I_0 \leftarrow \mathbf{w}_0$ 
2 Initialize an array  $\Delta G = []$ 
3 Initialize an array  $\Delta W = []$ 
4 Initialize  $last\_t = j_0$ 
5  $is\_explicit = False$ 
6 for  $t = 0; t < T; t++$  do
7   if  $(t - last\_t) \bmod T_0 == 0$  or  $t \leq j_0$  then
8     |  $is\_explicit = True$ 
9   else
10  end
11  if  $is\_explicit == True$  or  $t \leq j_0$  then
12    |  $last\_t = t$ 
13    | compute  $\nabla F(\mathbf{w}^I_t)$  exactly
14    | compute  $\nabla F(\mathbf{w}^I_t) - \nabla F(\mathbf{w}_t)$  based on the cached gradient  $\nabla F(\mathbf{w}_t)$ 
15    | /* check local convexity */
16    | if  $\langle \nabla F(\mathbf{w}^I_t) - \nabla F(\mathbf{w}_t), \mathbf{w}^I_t - \mathbf{w}_t \rangle \leq 0$  then
17    | | compute  $\mathbf{w}^I_{t+1}$  by using exact GD update (equation (1))
18    | | continue
19    | end
20    | set  $\Delta G[k] = \nabla F(\mathbf{w}^I_t) - \nabla F(\mathbf{w}_t)$ 
21    | set  $\Delta W[k] = \mathbf{w}^I_t - \mathbf{w}_t$ , based on the cached parameters  $\mathbf{w}_t$ 
22    |  $k \leftarrow k + 1$ 
23    | compute  $\mathbf{w}^I_{t+1}$  by using exact GD update (equation (1))
24  else
25    | Pass  $\Delta W[-m : ]$ ,  $\Delta G[-m : ]$ , the last  $m$  elements in  $\Delta W$  and  $\Delta G$ , which are from the  $j_1^{th}, j_2^{th}, \dots, j_m^{th}$ 
26    | iterations where  $j_1 < j_2 < \dots < j_m$  depend on  $t$ ,  $\mathbf{v} = \mathbf{w}^I_t - \mathbf{w}_t$ , and the history size  $m$ , to the
27    | L-BFGFS Algorithm (See Supplement) to get the approximation of  $\mathbf{H}(\mathbf{w}_t)\mathbf{v}$ , i.e.,  $\mathbf{B}_{j_m}\mathbf{v}$ 
28    | /* check local smoothness */
29    | if  $\|\mathbf{B}_{j_m}\mathbf{v}\| \geq \|\mathbf{v}\|$  then
30    | | go to line 12
31    | end
32    | Approximate  $\nabla F(\mathbf{w}^I_t) = \nabla F(\mathbf{w}_t) + \mathbf{B}_{j_m}(\mathbf{w}^I_t - \mathbf{w}_t)$ 
33    | Compute  $\mathbf{w}^I_{t+1}$  by using the "leave- $r$ -out" gradient formula, based on the approximated  $\nabla F(\mathbf{w}^I_t)$ 
34  end
35 end
36 return  $\mathbf{w}^I_t$ 
```

---



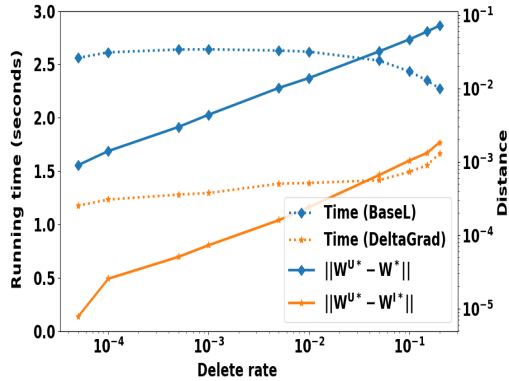


Figure 1: Running time and distance with varied deletion rate up to 20%

parameters (i.e.  $\mathbf{w}^{I^*}$  VS  $\mathbf{w}^{U^*}$ ) are still acceptable (on the order of  $10^{-3}$ ), far smaller than the error bound between  $\mathbf{w}^{U^*}$  and  $\mathbf{w}^*$  (on the order of  $10^{-1}$ ). Such a small difference between  $\mathbf{w}^{I^*}$  and  $\mathbf{w}^{U^*}$  also results in almost the same prediction performance, i.e.  $87.460 \pm 0.0011\%$  and  $87.458 \pm 0.0012\%$  respectively. This experiment thus provides some justification for the feasibility of DeltaGrad even when the number of the removed samples is not far smaller than the entire training dataset size.

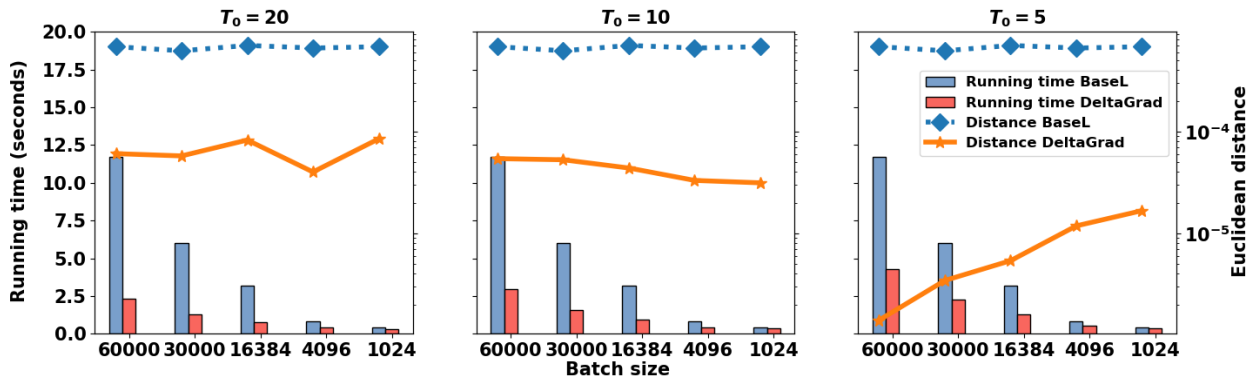


Figure 2: Running time and distance comparison with varying mini-batch size under fixed  $j_0 = 10$  and varying  $T_0$  ( $T_0 = 20$  VS  $T_0 = 10$  VS  $T_0 = 5$ )

## D.2 Influence of hyper-parameters on performance

To begin with, the influence of different hyper-parameters used in SGD and DeltaGrad is explored. We delete one sample from the training set of MNIST by running regularized logistic regression with the same learning rate and regularization rate as in Section 4 and varying mini-batch sizes (1024 - 60000),  $T_0$  ( $T_0 = 20, 10, 5$ ) and  $j_0$  ( $j_0 = 5, 10, 50$ ). The experimental results are presented in Figure 2-3. For different mini-batch sizes, we also used different epoch numbers to make sure that the total number of running iterations/steps in SGD are roughly the same. In what follows, we analyze how the mini-batch size, the hyper-parameters  $T_0$  and  $j_0$  influence the performance, thus providing some hints on how to choose proper hyper-parameters when DeltaGrad is used.

**Influence of the mini-batch size.** It is clear from Figure 2-3 that with larger mini-batch sizes, DeltaGrad can gain more speed with longer running time for both BaseL and DeltaGrad. As discussed in Section 4, to compute the gradients, other GPU-related overhead (the overhead to copy data from CPU

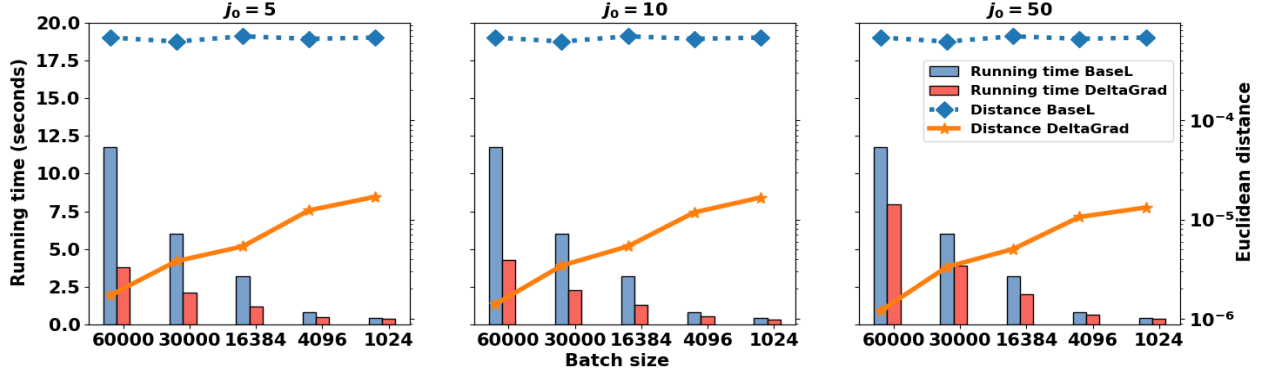


Figure 3: Running time and distance comparison with varying mini-batch size under fixed  $T_0 = 5$  and varying  $j_0$  ( $j_0 = 5$  VS  $j_0 = 10$  VS  $j_0 = 50$ )

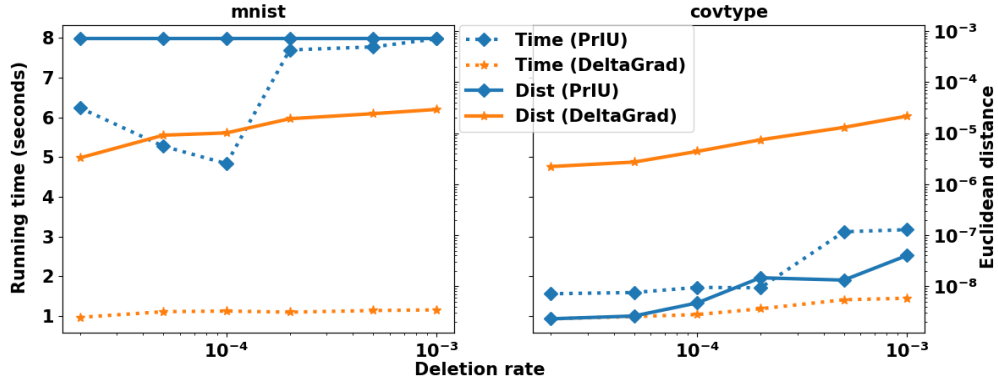


Figure 4: Comparison of DeltaGrad and PrIU

DRAM to GPU DRAM, the time to launch the kernel on GPU) cannot be ignored. This can become more significant when compared against the smaller computational overhead for smaller mini-batch data. Also notice that, when  $T_0 = 5$ , with increasing  $B$ , the difference between  $\mathbf{w}^U$  and  $\mathbf{w}^I$  becomes smaller and smaller, which matches our conclusion in Theorem 11, i.e. with larger  $B$ , the difference  $o(\frac{r}{n} + \frac{1}{B^{\frac{1}{4}}})$  is smaller.

**Influence of  $T_0$ .** By comparing the three sub-figures in Figure 2, the running time slightly (rather than significantly) decreases with increasing  $T_0$  for the same mini-batch size. This is explained by the earlier analysis in Section 4 on the non-ideal performance for GPU computation over small matrices. Interestingly, when  $T_0 = 10$  or  $T_0 = 20$ ,  $\|\mathbf{w}^{I,S} - \mathbf{w}^{U,S}\|$  does not decrease with larger mini-batch sizes. This is because in Formula (61), one component of the bound of  $\|\mathbf{w}^{I,S} - \mathbf{w}^{U,S}\|$  is

$$\frac{M_1(\frac{r}{n} + \frac{1}{B^{\frac{1}{4}}})}{C(1 - \frac{r}{n} - \frac{1}{B^{\frac{1}{4}}})} (1 - \eta C)^{yT_0} (1 - \eta\mu)^{j_0} d_{0,mT_0-1} \frac{1}{1 - (\frac{1-\eta\mu}{1-\eta C})^{T_0}}$$

(while the other component is  $o((\frac{r}{n} + \frac{1}{B^{\frac{1}{4}}}))$ ). Here  $d_{0,mT_0-1}$  increases with larger  $T_0$  and the term  $(1 - \eta C)^{yT_0}$  is not arbitrarily approaching 0 since  $yT_0$  cannot truly go to infinity. So when  $T_0 = 20$  and  $T_0 = 10$ , this component becomes the dominating term in the bound of  $\|\mathbf{w}^{I,S} - \mathbf{w}^{U,S}\|$ . So to make the bound  $o((\frac{r}{n} + \frac{1}{B^{\frac{1}{4}}}))$  hold, so that we can adjust the bound of  $\|\mathbf{w}^{I,S} - \mathbf{w}^{U,S}\|$  by varying  $B$ , proper choice of  $T_0$  is important. For example,  $T_0 = 5$  is a good choice for the MNIST dataset. This can achieve speed-ups comparable to larger  $T_0$  without sacrificing the closeness between  $\mathbf{w}^{I,S}$  and  $\mathbf{w}^{U,S}$ .

**Influence of  $j_0$ .** By comparing the three sub-figures in Figure 3, with increasing  $j_0$ , long “burn-in”

iterations are expected, thus incurring more running time. This, however, does not significantly reduce the distance between  $\mathbf{w}^{I,S}$  and  $\mathbf{w}^{U,S}$ . It indicates that we can select smaller  $j_0$ , e.g. 5 or 10 for more speed-up.

**Discussions on tuning the hyper-parameters for DeltaGrad.** Through our extensive experiments, we found that for regularized logistic regression, setting  $T_0$  as 5 and  $j_0$  as 5–20 would lead to some of the most favorable trade-offs between running time and the error  $\|\mathbf{w}^{U,S} - \mathbf{w}^{I,S}\|$ . But in terms of more complicated models, e.g. 2-layer DNN, higher  $j_0$  (even half of the total iteration number) and smaller  $T_0$  (2 or 1) are necessary. Similar experiments were also conducted on adding training samples, in which similar trends were observed.

### D.3 Comparison against the state-of-the-art work

To our knowledge, the closest work to ours is Wu et al. (2020), which targets simple ML models, i.e. linear regression and regularized logistic regression with an ad-hoc solution (called PrIU) rather than solutions for general models. Their solutions can only deal with the deletion of samples from the training set without supporting the addition of samples. In our experiments, we compared DeltaGrad (with  $T_0 = 5$  and  $j_0 = 10$ ) against PrIU by running regularized logistic regression over MNIST and covtype with the same mini-batch size (16384), the same learning rate and regularization rate, but with varying deletion rates.

**Table 1: Memory usage of DeltaGrad and PrIU(GB)**

| Deletion rate      | MNIST |           | covtype |           |
|--------------------|-------|-----------|---------|-----------|
|                    | PrIU  | DeltaGrad | PrIU    | DeltaGrad |
| $2 \times 10^{-5}$ | 26.61 | 2.74      | 9.30    | 2.56      |
| $5 \times 10^{-5}$ | 27.02 | 2.74      | 9.30    | 2.56      |
| $1 \times 10^{-4}$ | 27.13 | 2.74      | 9.30    | 2.55      |
| $2 \times 10^{-4}$ | 27.75 | 2.74      | 9.31    | 2.56      |
| $5 \times 10^{-4}$ | 29.10 | 2.74      | 10.67   | 2.56      |
| $1 \times 10^{-3}$ | 29.10 | 2.74      | 10.67   | 2.56      |

The running time and the distance term  $\|\mathbf{w}^U - \mathbf{w}^I\|$  of both PrIU and DeltaGrad with varying deletion rate are presented in Figure 4. First, it shows that DeltaGrad is always faster than PrIU, with more significant speed-ups on MNIST. The reason is that the time complexity of PrIU is  $O(rp)$  for each iteration where  $p$  represents the total number of model parameters while  $r$  represents the reduced dimension after Singular Value Decomposition is conducted over some  $p \times p$  matrix. This is a large integer for large sparse matrices, e.g. MNIST.

As a result,  $O(rp)$  is larger than the time complexity of DeltaGrad. Also, the memory usage of PrIU and DeltaGrad is shown in Table 1. PrIU needs much more DRAM (even 10x in MNIST) than DeltaGrad. The reason is that to prepare for the model update phase, PrIU needs to collect more information during the training phase over the full dataset. This is needed in the model update phase and is quadratic in the number of the model parameters  $p$ . The authors of Wu et al. (2020) claimed that their solution cannot provide good performance over sparse datasets in terms of running time, error term  $\mathbf{w}^U - \mathbf{w}^I$  and memory usage. In contrast, both the time and space overhead of DeltaGrad are smaller, which thus indicates the potential of its usage in the realistic, large-scale scenarios.

### D.4 Experiments on large ML models

In this section, we compare DeltaGrad with BaseL using the state-of-the-art ResNet152 network (He et al., 2016) (ResNet for short hereafter) with all but the top layer frozen, for which we use the pre-trained parameters from Pytorch torchvision library<sup>1</sup>. The pre-trained layers with fixed parameters are regarded as the feature transformation layer, applied over each training sample as the pre-processing step before the training phase. Those transformed features are then used to train the last layer of ResNet, which is thus equivalent to training a logistic regression model.

This experiment is conducted on CIFAR-10 dataset (Krizhevsky et al., 2009), which is composed of 60000  $32 \times 32$  color images (50000 of them are training samples while the rest of are test samples). We run

<sup>1</sup><https://pytorch.org/docs/stable/torchvision/models.html>

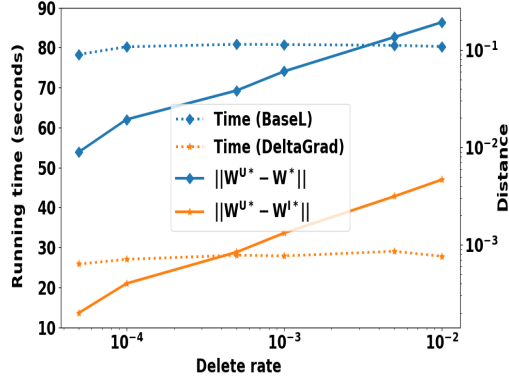


Figure 5: Comparison of DeltaGrad and BaseL on the CIFAR-10 dataset with pre-trained ResNet152 network

SGD with mini-batch size 10000, fixed learning rate 0.05 and L2 regularization rate 0.0001. Similar to the experimental setup introduced in Section 4 in the main paper, the deletion rate is varied from 0 to 1% and the model parameters are updated by using BaseL and DeltaGrad (with  $T_0$  as 5 and  $j_0$  as 20) respectively after the deletion operations. The experimental results are presented in Figure 5, again showing significant speed-ups for DeltaGrad relative to BaseL (up to 3x speed-ups when the deletion rate is 0.005%) with far smaller error bound (up to  $4 \times 10^{-3}$ ) than the baseline error bound (up to  $2 \times 10^{-2}$ ). Since it is quite common to reuse sophisticated pre-trained models in practice, we expect that the use of DeltaGrad in this manner is applicable in many cases.

## D.5 Applications of DeltaGrad to robust learning

As Section 5 in the main paper reveals, DeltaGrad has many potential applications. In this section, we explored how DeltaGrad can accelerate the evaluations of the effect of the outliers in robust statistical learning. Here the effect of outliers is represented by the difference of the model parameters before and after the deletion of the outliers (see Yu and Yao (2017)).

In the experiments, we start by training a model on the training dataset (RCV1 here) along with some randomly generated outliers. Then we remove those outliers and update the model on the remaining training samples by using DeltaGrad and BaseL. We also evaluate the effect of the fraction of outlier: the ratio between the number of the outliers and the training dataset size is also defined as the *Deletion rate*. It is varied from 1% to 10%. According to the experimental results shown in Figure 6, when there are up to 10% outliers in the training dataset, DeltaGrad is at least 2.18x faster than BaseL in evaluating the updated model parameters by only sacrificing little computational accuracy (no more than  $5 \times 10^{-3}$ ), thus reducing the computational overhead on evaluating the effect of the outliers in robust learning.

## References

- L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- R. H. Byrd, J. Nocedal, and R. B. Schnabel. Representations of quasi-newton matrices and their use in limited memory methods. *Mathematical Programming*, 63(1-3):129–156, 1994.
- K. Chaudhuri and C. Monteleoni. Privacy-preserving logistic regression. In *Advances in neural information processing systems*, pages 289–296, 2009.
- A. R. Conn, N. I. Gould, and P. L. Toint. Convergence of quasi-newton matrices generated by the symmetric rank one update. *Mathematical programming*, 50(1-3):177–195, 1991.
- C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- A. Ginart, M. Guan, G. Valiant, and J. Y. Zou. Making ai forget you: Data deletion in machine learning. In *Advances in Neural Information Processing Systems*, pages 3513–3526, 2019.

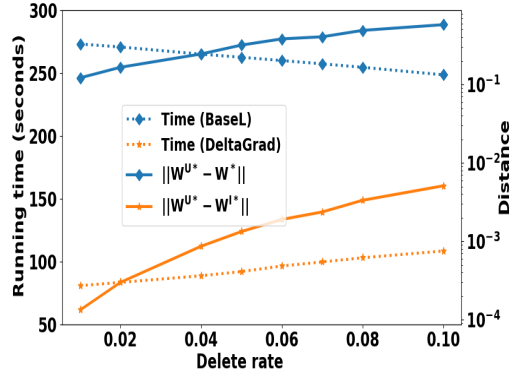


Figure 6: Comparison of DeltaGrad and BaseL on RCV1 dataset after deleting outliers

- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- R. I. Oliveira. Concentration of the adjacency matrix and of the laplacian in random graphs with independent edges. *arXiv preprint arXiv:0911.0600*, 2009.
- J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- J. A. Tropp. The expected norm of a sum of independent random matrices: An elementary approach. In *High dimensional probability VII*, pages 173–202. Springer, 2016.
- Y. Wu, V. Tannen, and S. B. Davidson. Priu: A provenance-based approach for incrementally updating regression models. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 447–462, 2020.
- C. Yu and W. Yao. Robust linear regression: A review and comparison. *Communications in Statistics-Simulation and Computation*, 46(8):6261–6282, 2017.